

A Bayesian approach for analysis of whole-genome bisulphite sequencing data identifies disease-associated changes in DNA methylation

Running title: ABBA for analysis of WGBS in disease

Owen J.L. Rackham^{1,\$}, Sarah R. Langley^{1,\$}, Thomas Oates^{2,\$}, Eleni Vradi^{3,\$}, Nathan Harmston¹, Prashant K. Srivastava⁴, Jacques Behmoaras⁵, Petros Dellaportas^{6,7,*}, Leonardo Bottolo^{8,7,9,*} & Enrico Petretto^{1,2*}

¹Duke-NUS Medical School, Singapore. ²MRC London Institute of Medical Sciences (LMS), Imperial College London, UK. ³Department of Statistics, Athens University of Economics and Business, Athens, GR. ⁴Division of Brain Sciences, Faculty of Medicine, Imperial College London, UK. ⁵Centre for Complement and Inflammation Research, Imperial College London, UK. ⁶Department of Statistical Science, University College London, UK. ⁷The Alan Turing Institute, London, UK. ⁸Department of Medical Genetics, University of Cambridge, UK. ⁹MRC Biostatistics Unit, Cambridge Institute of Public Health, Cambridge, UK. \$ These authors contributed equally to this work.

*Corresponding authors:

Department of Statistical Science, University College, Gower Street, London, WC1E 6BT E-mail: p.dellaportas@ucl.ac.uk; and Department of Medical Genetics, University of Cambridge, Box 238, Lv 6 Addenbrooke's Treatment Centre, Addenbrooke's Hospital, Cambridge Biomedical Campus, Cambridge, CB2 0QQ, UK. E-mail: lb664@cam.ac.uk; and Duke-NUS Medical School, 8 College road 169857, Republic of Singapore. E-mail: enrico.petretto@duke-nus.edu.sg

1 ABSTRACT

2 DNA methylation is a key epigenetic modification involved in gene regulation whose contribution to
3 disease susceptibility remains to be fully understood. Here, we present a novel Bayesian smoothing
4 approach (called ABBA) to detect differentially methylated regions (DMRs) from whole-genome bisulphite
5 sequencing (WGBS). We also show how this approach can be leveraged to identify disease-associated
6 changes in DNA methylation, suggesting mechanisms through which these alterations might affect disease.
7 From a data modeling perspective, ABBA has the distinctive feature of automatically adapting to different
8 correlation structures in CpG methylation levels across the genome whilst taking into account the distance
9 between CpG sites as a covariate. Our simulation study shows that ABBA has greater power to detect
10 DMRs than existing methods, providing an accurate identification of DMRs in the large majority of
11 simulated cases. To empirically demonstrate the method's efficacy in generating biological hypotheses, we
12 performed WGBS of primary macrophages derived from an experimental rat system of
13 glomerulonephritis and used ABBA to identify >1,000 disease-associated DMRs. Investigation of these
14 DMRs revealed differential DNA methylation localized to a 600bp region in the promoter of the *Ifitm3*
15 gene. This was confirmed by ChIP-seq and RNA-seq analyses, showing differential transcription factor
16 binding at the *Ifitm3* promoter by JunD (an established determinant of glomerulonephritis) and a
17 consistent change in *Ifitm3* expression. Our ABBA analysis allowed us to propose a new role for *Ifitm3* in
18 the pathogenesis of glomerulonephritis via a mechanism involving promoter hypermethylation that is
19 associated with *Ifitm3* repression in the rat strain susceptible to glomerulonephritis.

20

1 INTRODUCTION

2 One of the most important epigenetic modifications directly affecting DNA is methylation, where a
3 methyl group is added to a cytosine base in the DNA sequence creating 5-methylcytosine. High-
4 throughput sequencing techniques, such as whole-genome bisulphite sequencing (WGBS), now allow for
5 genome-wide methylome data to be collected at single base-resolution (Harris *et al.* 2010). However, the
6 challenge remains on how to accurately identify DNA methylation changes at the genome-wide level and
7 also account for the complex correlation structures present in the data. Whilst it is still not fully
8 understood how DNA methylation affects gene expression, it has been shown that depending on the
9 location of the modification it can either have a positive or negative effect on the level of expression of
10 genes (Gutierrez-Arcelus *et al.* 2013). How methylation patterns are regulated is complex and a full
11 understanding of this process requires elucidating the mechanisms for *de novo* DNA methylation and
12 demethylation, as well as the maintenance of methylation (Chen and Riggs 2011) However, the majority of
13 functional methylation changes are found in methylation sites where cytosines are immediately followed
14 by guanines, known as CpG dinucleotides (Ziller *et al.* 2011). These are not positioned randomly across
15 the genome but tend to appear in clusters called CpG islands (CpGI) (Deaton and Bird 2011). It has been
16 also shown that there are concordant methylation changes within CpGI and in the genomic regions
17 immediately surrounding CpGI (also known as CpGI shores or CpGS). These “spatially correlated” DNA
18 methylation patterns tend to be more strongly associated with gene expression changes than the
19 methylation changes occurring in other parts of the genome (Gutierrez-Arcelus *et al.* 2015). The
20 correlation of methylation levels between CpG sites is also highly dependent on their genomic context,
21 varying greatly depending on where in the genome they are located (Zhang *et al.* 2015). For
22 computational convenience, the dependence of methylation patterns between CpG sites is sometimes
23 ignored by methods for differential methylation analysis. Alternatively, a simplified estimation of the
24 correlation of methylation levels between neighbouring CpG sites (Bell *et al.* 2011) based on a user-
25 defined parameterization of the degree of smoothing is introduced. These strategies might not be
26 appropriate across different experimental scenarios and instead we propose an automatic probabilistic
27 smoothing procedure of the average methylation levels across replicates (hereafter methylation profiles).

28 Beyond the initial univariate analysis of methylation changes at each individual CpG (for instance,
29 using the Fisher’s exact test), recently the focus has shifted to identifying differentially methylated regions
30 (DMRs), since coordinated changes in CpG methylation across genomic regions are known to impart the
31 strongest regulatory influence. To this aim, a number of tools have been proposed to detect DMRs from

1 WGBS data. Typically, these methods normally take one of two approaches: Either model the number of
2 methylated/unmethylated reads using a binomial, negative-binomial distribution or discrete distributions
3 with an over-dispersion parameter) such as MethylKit (Akalın *et al.* 2012), MethylSig (Park *et al.* 2014)
4 and DSS (Feng *et al.* 2014). Alternatively in order to account for the correlation of methylation profiles
5 between neighbouring CpG sites, a smoothing operator is applied in tools like BSmooth (Hansen *et al.*
6 2012), BiSeq (Hebestreit *et al.* 2013), DSS-single (Wu *et al.* 2015) – reviewed in Robinson *et al.* 2014 and
7 Yu and Sun 2016b). Methods based on spline- (Hansen *et al.* 2012), kernel- (Hebestreit *et al.* 2013)
8 perform generally well in practical applications. However, their results and the identification of the DMRs
9 depend on the choice of the smoothing parameters values, e.g., window size or kernel bandwidth, a
10 feature that makes them less general and prone to perform unequally when the default parameters values
11 are changed. In these cases, smoothing parameters tuned by time-consuming sensitivity analysis based on
12 different parameterizations is usually recommended, although this strategy is rarely applied in real data
13 analyses. Other approaches, e.g., metilene (Jühling *et al.* 2015), propose segmentation algorithms to detect
14 DMRs between single/groups of replicates without making any model assumption about the data
15 generating mechanism and less dependent on the parameters definition. Furthermore, several other
16 algorithms have been introduced, e.g., MOABS (Sun *et al.* 2014), Lux (Äijö *et al.* 2016), and MACAU (Lea *et*
17 *al.* 2015), showing that bisulfite sequencing data analysis is an active area of research.

18 To address this dependence on parameterization and the subsequent lack of generality, we propose a
19 fully Bayesian approach, approximate Bayesian bisulphite sequencing analysis (or ABBA), designed to
20 smooth automatically the underlying - not directly observable - methylation profiles and reliably identify
21 DMRs whilst borrowing information vertically across biological replicates and horizontally across
22 correlated CpGs (**Fig. 1**). We highlight that this fully Bayesian specification is not adopted by previous
23 DMR detection techniques, owing to the computational overhead of the inferential procedure. We address
24 the high computational demands by utilizing a highly efficient inferential tool (Rue *et al.* 2009) for
25 Bayesian models (see below and **Methods**). To demonstrate the benefits of adopting ABBA over existing
26 approaches, we report a comprehensive simulation study where we benchmarked ABBA against five
27 commonly used alternative methods (Fisher Exact Test, BSmooth, MethylKit, MethylSig, DSS) and
28 considered a proposed new one (metilene) and assessed the effect of a different biological and
29 experimental conditions (by varying parameters related to data integrity and quality of the signal) on the
30 performance of each method. The results from this benchmark clearly indicate that ABBA is the best
31 performing method, being both robust to changes in factors affecting data quality (e.g., sequencing
32 coverage, errors associated with the methylation call) and level of noise in methylation signal. To

1 benchmark our proposed method on a real dataset, we generated new WGBS data in macrophages from
2 an established rat model of glomerulonephritis (Aitman *et al.* 2006) and control strain, and used ABBA for
3 the genome-wide identification of DMRs. An additional comparison performed with the best alternative
4 method (that arose from the simulation study) showed that ABBA has increased power to detect changes
5 in DNA methylation involving genes and pathways relevant to glomerulonephritis. Furthermore, this
6 comparison exemplifies how the DMR results obtained by alternative approaches depend heavily on the
7 choice of relevant smoothing parameters (e.g., window size used in DSS). We also integrated the DMR
8 results of ABBA with transcription factor binding site analysis, RNA-seq and ChIP-seq data generated in
9 the same system, and in this we revealed a previously unappreciated role for the *Ifitm3* gene in the
10 pathogenesis of glomerulonephritis, providing a proof of concept for real data applications of the ABBA
11 approach.

12

13 MATERIALS AND METHODS

14 Below we reported the key aspects of the latent Gaussian model and Integrated Nested Laplace
15 Approximation (INLA). Interested reader can also refer to Rue and Martino 2007 and Rue *et al.* 2009.

16 **Latent Gaussian model.** A latent Gaussian model (LGM) can be described by a three-stage hierarchical
17 model

$$18 \quad y_i | x_i, \boldsymbol{\theta} \sim \pi(y_i | x_i, \boldsymbol{\theta}), \quad (1)$$

$$19 \quad \mathbf{x} | \boldsymbol{\theta} \sim N(\boldsymbol{\mu}(\boldsymbol{\theta}), \mathbf{Q}^{-1}(\boldsymbol{\theta})), \quad (2)$$

$$20 \quad \boldsymbol{\theta} \sim \pi(\boldsymbol{\theta}), \quad (3)$$

21 where $y_i, i = 1, \dots, n$, are the observed values, \mathbf{x} is n -dimensional vector of latent variables and $\boldsymbol{\theta}$ is p -
22 dimensional vector of model parameters. (1) is the *observations equation* and it describes the probabilistic
23 model for each observation conditionally on the latent variable x_i and the model parameters $\boldsymbol{\theta}$, (2) is the
24 *latent Gaussian field equation* with the latent variables distributed as a p -dimensional normal distribution,
25 with mean vector $\boldsymbol{\mu}(\boldsymbol{\theta})$ and a sparse precision matrix $\mathbf{Q}(\boldsymbol{\theta})$. Both quantities can depend on the model
26 parameters vector $\boldsymbol{\theta}$ whose distribution is described in the *parameter equation* (3). The Gaussian vector \mathbf{x}
27 exhibits a particular conditional dependence (or Markov) structure which is reflected in its precision
28 matrix $\mathbf{Q}(\boldsymbol{\theta})$.

1 **Integrated Nested Laplace Approximation.** INLA is a computational approach to perform statistical
 2 inference for LGM. It provides a fast and accurate alternative to exact MCMC (Gilks *et al.* 1996) and other
 3 sampling-based methods such as Sequential Monte Carlo (Doucet *et al.* 2001). They become prohibitively
 4 computationally expensive when the length of the sequence considered is too long, resulting in infeasible
 5 run times. The INLA solution with a mix of Laplace approximations (Tierney and Kadane 2012) and
 6 numerical integrations offers a pragmatic inferential tool to fit LGMs and it provides answers in hours
 7 whereas MCMC requires days. The INLA inferential procedure consists of three steps:

- 8 1. Compute the approximation to the marginal posterior $\pi(\boldsymbol{\theta}|\mathbf{y})$ and by-product to $\pi(\theta_j|\mathbf{y}), j = 1, \dots, p$;
- 9 2. Compute the approximation to $\pi(x_i|\mathbf{y}, \boldsymbol{\theta}), i = 1, \dots, n$;
- 10 3. Combine 1 and 2 above and compute the approximation to the marginal posterior $\pi(x_i|\mathbf{y})$.

11 **ABBA model.** Based on LGM, the ABBA model can be described by a three-stage hierarchical model:

$$12 \quad y_{igr} | \pi_{igr} \sim \text{Binomial}(n_{igr}, \pi_{igr}), \quad (4)$$

$$13 \quad \text{logit}(\pi_{igr}) | \sigma_g^2 \sim \text{N}(\mu_{ig}, \sigma_g^2) \quad (5)$$

$$14 \quad \mu_{ig} | \rho_{ig}^2 \sim \text{N}(\mu_{i-1,g}, \rho_{ig}^2) \quad (6)$$

$$15 \quad \sigma_g^{-2} \sim \text{Gam}(0.1, 0.1) \quad (7)$$

$$16 \quad \rho_{ig}^2 = \rho_g^2 | p_i - p_{i-1} | \quad \text{with} \quad \rho_g^{-2} \sim \text{Gam}(0.1, 0.1) \quad (8)$$

17 (4) is the *first part of the observations equation* where $i = 1, \dots, m$ denotes the CpG, $g = 1, 2$ the group (e.g.,
 18 case and control group), and $r = 1, \dots, R$ the experimental replicate. y_{igr}, n_{igr} and π_{igr} are the observed
 19 number of methylated reads, the read depth and the proportion of methylation for the i th CpG site, g th
 20 group and r th experimental replicate, respectively. (5) is the *second part of the observations equation* and
 21 it describes a random effect across the experimental replicates with a specific variance σ_g^2 for each group.
 22 In (5), $\text{logit}(z)$ indicates the logit transformation, $\text{logit}(z) = \log(1/(1-z))$. The observation equation (5)
 23 assumes that the methylation proportions are drawn from the same distribution within each group but
 24 are different between groups.

25 (6) is the *latent Gaussian field (LGF) equation*. The dependence of the DNA methylation pattern between
 26 CpGs is modelled as a non-stationary random walk of order 1, RW(1): μ_{ig} follows a normal distribution
 27 with mean $\mu_{i-1,g}$ (defined in the $(i-1)$ th CpG) and variance ρ_{ig}^2 which is specific for each CpG and group.
 28 (5) and (6) highlight an important feature of ABBA model that is able to model vertically the information

1 contained in the replicates by a random effect model and horizontally the information about the CpG
 2 methylation levels correlation by a LGF.

3 The model is completed by specification in (7) and (8) of the random effect and LGM prior precision, i.e.
 4 the inverse of the variance. For computational convenience we introduce a CpG site spacing and
 5 decompose ρ_{ig}^2 into $\rho_g^2 |p_i - p_{i-1}|$, where ρ_g^2 is the global smoothing parameter specific for each group that
 6 needs to be estimated and p_i and p_{i-1} are the chromosomal locations of two consecutive CpG sites. This
 7 implies that the correlation between μ_{ig} and $\mu_{i-1,g}$, depends on the distance between the two consecutive
 8 CpG sites and in particular, it decreases as this distance increases in keeping with empirical evidence
 9 reported in Bell JT et al. 2011, Zhang (2015) and in our real data set (see **Supplementary Figure 1**). With
 10 this formulation only σ_g^2 and ρ_g^2 need to be estimated for each group. It also implies a sparse precision
 11 matrix $\mathbf{Q}(\boldsymbol{\theta})$ for the LGF in (2) making the overall inferential process efficient.

12 Finally, non-informative priors are assigned to the precision parameters σ_g^{-2} and ρ_g^{-2} which are
 13 distributed as a gamma density with mean 1 and variance 10 (default INLA values). Sensitivity analysis on
 14 the gamma density parameterization shows no departure from the results obtained using the default
 15 values. See **Supplementary Table 1** for details on the posterior density of σ_g^{-2} and ρ_g^{-2} under INLA
 16 default and alternative parameterization on selected simulated examples.

17 When a single replicate is available, since $\sigma_g^2 = 0$, (4) and (5) simplify to

$$18 \quad y_{ig} | \pi_{ig} \sim \text{Binomial}(n_{ig}, \pi_{ig}),$$

$$19 \quad \text{logit}(\pi_{ig}) = \mu_{ig}.$$

20 While some methods for DMR detection (Feng *et al.* 2014; Wu *et al.* 2015), allow for over-dispersion by
 21 assuming a beta-binomial model, (4) and (5) imply a logistic-normal model. After integrating out (6),
 22 $\int N(\text{logit}(\pi_{igr}) | \mu_{ig}, \sigma_g^2) N(\mu_{ig} | \mu_{i-1,g}, \rho_{ig}^2) d\mu_{ig} = N(\text{logit}(\pi_{igr}) | \mu_{i-1,g}, \sigma_g^2 + \rho_{ig}^2)$, it can be shown that
 23 marginally

$$24 \quad V(Y_{ig} | \sigma_g^2, \rho_{ig}^2) \approx n_{ig} \pi_{ig} (1 - \pi_{ig}) \{1 + (\sigma_g^2 + \rho_{ig}^2) (n_{ig} - 1) \pi_{ig} (1 - \pi_{ig})\},$$

25 where $\pi_{ig} \equiv \exp(\mu_{ig}) / (1 + \exp(\mu_{ig}))$. The above equation illustrates that *a priori* the marginal degree of
 26 variability per CpG site under ABBA model is the variance of the binomial model multiplied by an over-
 27 dispersion factor that depends on the combined effect of σ_g^2 , the replicates variability, and ρ_{ig}^2 , the

1 variance of the unobserved methylation profile. When a single replicate is available, the over-dispersion
2 depends only on ρ_{ig}^2 .

3 **ABBA algorithm.** The ABBA algorithm consists of two steps:

4 1. Compute the approximation to the marginal posteriors of σ_g^2 , the variance of the random effect, and ρ_g^2 ,
5 $g = 1, 2$ the smoothing parameters; given the model specification $\rho_{ig}^2 = \rho_g^2 |p_i - p_{i-1}|$, it is also possible
6 to derive the marginal posteriors of ρ_{ig}^2 ;

7 2. Compute the approximation to marginal posterior $\pi(\mu_{ig}|\mathbf{y})$, where $\mathbf{y} = (y_{igr})_{i=1, \dots, n; g=1, 2; r=1, \dots, R}$; then
8 the marginal posterior of the unobserved methylation profile $\pi(\pi_{ig}|\mathbf{y})$ is obtained by using the inverse
9 logit transformation of μ_{ig} , $z \equiv \exp\{\text{logit}(z)\} / [1 + \exp\{\text{logit}(z)\}]$.

10 **Global differential methylation and FDR calculation.** ABBA inference about DMRs is based on the
11 posterior methylation probability (PMP) $\pi(\pi_{ig}|\mathbf{y})$ and the posterior differential methylation probability
12 (PDMP) $\pi(\pi_{i1}|\mathbf{y}) - \pi(\pi_{i2}|\mathbf{y})$. The posterior mean methylation probability $E(\pi_{ig}|\mathbf{y})$ summarizes the
13 information contained in the PMP and it is used to define the posterior mean differential methylation
14 between two groups, $d_i = E(\pi_{i1}|\mathbf{y}) - E(\pi_{i2}|\mathbf{y})$. Once the LGF has been integrated out by INLA inferential
15 process, $\pi(\pi_{ig}|\mathbf{y})$, $i = 1, \dots, n$, and in turn d_i s become marginally independent. This allows the
16 straightforward application of a non-parametric false discovery rate (FDR) procedure without the burden
17 of correlated signals. To distinguish between the null distribution (no differential methylation) and the
18 alternatives, we fit a mixture of three truncated normal densities

$$19 \quad d_i \sim \pi_- N_{[-1,1]}(\theta_-, \xi_-^2) + \pi_0 N_{[-1,1]}(\theta_0, \xi_0^2) + \pi_+ N_{[-1,1]}(\theta_+, \xi_+^2), \quad (9)$$

20 where $N_{[-1,1]}$ is a normal density truncated between $[-1,1]$, $\pi_-, \pi_0, \pi_+ \in (0,1)$ with $\pi_- + \pi_0 + \pi_+ = 1$ are
21 the mixing weights of the “negative” differentially methylated, no differentially methylated and “positive”
22 differentially methylated with respect the control group, respectively, $\theta_-, \theta_0, \theta_+$ are the unknown centers
23 of the differentially methylated groups and $\xi_-^2, \xi_0^2, \xi_+^2$ are the unknown variances. Under the null
24 hypothesis we set $\theta_0 = 0$. For identifying the components of mixture model we also impose the condition
25 $\pi_0 \geq \pi_- + \pi_+$ under the assumption that the large majority of CpG sites are not differentially methylated.

26 Although the choice of a three component mixture model works well in real data examples (see
27 **Supplementary Figure 2**), this assumption can be relaxed. For instance, as suggested in Sun and Cai
28 2009, the non-null distribution f_1 can have more than two components. This allows a better fitting of the
29 tails of distribution of d_i 's and the identification of more than two differentially methylated groups. For

1 instance the choice of the number of components can be based on Bayesian Information Criterion (BIC).
 2 However this requires running the FDR procedure several times for each choice of the number of
 3 components. Another possibility which is less computational intensive relies on the approximation of f_1
 4 by using a non-parametric Gaussian kernel density estimation (Kuan and Chiang 2012).

5 Maximum likelihood estimates of (9) are obtained by the EM algorithm (Dempster *et al.* 1977) taking
 6 particular care to avoid local maxima in the likelihood surface by running the EM algorithm from different
 7 starting points. Using the EM algorithm, the posterior probability of a CpG site belonging to each of the
 8 three component is

$$9 \quad P(z_i = "-") = \frac{\pi_- N_{[-1,1]}(d_i; \theta_-, \xi_-^2)}{C},$$

$$10 \quad P(z_i = "0") = \frac{\pi_0 N_{[-1,1]}(d_i; 0, \xi_0^2)}{C},$$

$$11 \quad P(z_i = "+") = \frac{\pi_+ N_{[-1,1]}(d_i; \theta_+, \xi_+^2)}{C}$$

12 with $C = \pi_- N_{[-1,1]}(d_i; \theta_-, \xi_-^2) + \pi_0 N_{[-1,1]}(d_i; 0, \xi_0^2) + \pi_+ N_{[-1,1]}(d_i; \theta_+, \xi_+^2)$.

13 Similarly to Broët *et al.* 2004, for a constant t , we define the estimated $\widehat{\text{FDR}}(t)$ as

$$14 \quad \widehat{\text{FDR}}(t) = \frac{\sum_{i \in \mathcal{J}_-} P(z_i = "0") + \sum_{i \in \mathcal{J}_+} P(z_i = "0")}{n_- + n_+} \quad (10)$$

15 where $\mathcal{J}_- = \{i: d_i \leq -t\}$, $\mathcal{J}_+ = \{i: d_i \geq t\}$, $n_- = \#(\mathcal{J}_-)$ and $n_+ = \#(\mathcal{J}_+)$. (10) defines the global FDR as the
 16 average local FDR which, for posterior probabilities, is defined as $1 - P(z_i = "-") - P(z_i = "+") =$
 17 $P(z_i = "0")$. Finally the constant t is chosen such that $\widehat{\text{FDR}}(t) \leq \text{FDR}$.

18 In summary, the FDR procedure for ABBA consists of two steps:

- 19 1. Fit a mixture of truncated normal densities with three components on the d_i s values; obtain the
 20 posterior probability that each d_i belongs to each of the three components;
- 21 2. Calculate the constant t such that $\widehat{\text{FDR}}(t) \leq \text{FDR}$ for a desired level of FDR;

22 For computational efficiency our FDR procedure can be run on each chromosome separately and then the
 23 results can be aggregated at the genome-wide level (Efron 2008). Besides the computational speed, this
 24 strategy does not assume the existence of a global methylation level difference between the two
 25 conditions that may not hold in practice. The separate-class model (Efron 2008), can be used to combine
 26 separate chromosome-wide FDRs.

1 **WGBS data simulation.** WGBS data have a number of intrinsic characteristics that can vary depending on
2 the cell-types/tissue complexity being studied or on technical issues related to the sequencing. In order
3 assess which method is the most robust for analyzing WGBS data it is important that changes in each of
4 these characteristics are taken into account. Here we take advantage of our previously published WGBS-
5 data simulator (Rackham *et al.* 2015) that allows us to generate unbiased benchmarking datasets with
6 several varying parameters. Wherever possible we will refer to the notation used in Rackham *et al.* 2015;
7 the parameters are the following:

- 8 1. Number of replicates - the parameter r was set to vary between $r = 1,2,3$ within each group;
- 9 2. Average read depth - at each CpG site for all replicates and groups, the number of reads n_{igr} ,
10 $i = 1, \dots, m$ and $g = 1,2$, is simulated using a Poisson distribution with average read depth λ . The
11 parameter λ was set to be either 10 or 30 reads on average per CpG site;
- 12 3. Level of noise - the parameter s_0 controls the level of noise added the probability of methylation at
13 each CpG site for all replicates and groups and simulates the measurement error resulting from the
14 sampling of DNA segments during sequencing

$$15 \quad \pi_{irg} = \text{logit}^{-1}(\text{logit}(\pi_{rg}) + \varepsilon_i),$$

16 where π_{rg} is the global probability of methylation of the binomial (emission) distribution based on
17 the real dataset analyzed (see details in Rackham *et al.* 2015) and $\varepsilon_i \sim N(0, s_0)$, $i = 1, \dots, m$. s_0 was set
18 to vary between 0.1, 0.2 and 0.3 to model different level of noise. To calibrate the value of s_0 ,
19 **Supplementary Table 2** provides a Monte Carlo estimation of the effect of different values of the
20 noise level on π_{irg} .

- 21 4. Methylation probability difference - the parameter $\Delta meth$ reported in Rackham *et al.* 2015 as “phase
22 difference” controls the magnitude of the difference between the probabilities of methylation in each
23 group and was set to vary between 20%, 30%, 50% or 70%. This difference is obtained on CpG sites
24 where both case and control samples share the same methylated status (methylated or
25 unmethylated), by adding a given value to the probability in either cases or controls. The total length
26 of the sequence where this difference appears is no greater than 5% (WGBSSuite default value) of the
27 total length of the simulated region.
- 28 5. We also considered an additional parameter δ (not available for modeling in WGBSSuite), which
29 introduces a further error associated with the methylation call. After selecting at random with a given
30 probability δ a CpG site in the g th group for all replicates, we switch its methylation status between
31 the two groups. In our simulation study, the parameter δ has been varied from 0, 0.05 and 0.1.

1 To perform the benchmarking we generate 5 replicates of 5,000 CpGs for each combination of the
2 above parameters. The resulted in a total of 216 benchmarking datasets (3 cases for the number of
3 replicates, 2 cases for the average read depth, 3 cases for the level of noise, 4 cases for the methylation
4 probability difference, 3 cases for the parameter δ) which are replicated 5 times (5,400,000 CpGs in total)
5 to assess the Monte Carlo average performance for each combination of parameters. In these datasets the
6 size of the differentially methylated regions has a median size of 15 CpGs (see **Supplementary Figure 3**).
7 The proportion of differentially methylated CpGs cannot exceed 20% of all CpGs (i.e., ~1000 CpGs).

8 **Receiver operator curve (ROC) construction for benchmarking.** In order to generate the ROC curve
9 the performance is calculated CpG-wise. For a given DMR, detection of each of the CpG contained within is
10 considered as a true positive, whilst CpGs that are not detected are considered false negatives. Outside of
11 the DMR the opposite criteria is applied. We choose this assignment criteria rather than calling detection
12 of a each DMR since it provides a useful quantification of the extent each DMR is captured by each
13 technique, for instance if one technique correctly identifies all the CpGs in a DMR, the method is deemed to
14 perform better than an approach that identifies correctly only 80% of the CpGs within the same DMR.

15 **WGBS data pre-processing for ABBA.** To run ABBA efficiently at the genome-wide level we took
16 advantage of cluster-computing environment that enables parallel computation, and to this aim we
17 preprocessed the WGBS data as follows. After the raw WGBS data were aligned, we removed CpG sites
18 where less than 50% of the samples contain reads. Next, we split the WGBS data into chunks such that the
19 distance between the last CpG site in one chunk and the first CpG in the next chunk is greater than
20 3,000bp. It has been previously shown that the correlation of DNA methylation levels between CpG sites
21 decreases dramatically after 400bp (Zhang *et al.* 2015), so splitting the data in this way implies a
22 particular conditional dependence structure in our data defined by a sparse block-diagonal precision
23 matrix $\mathbf{Q}(\theta)$ where each block corresponds to a WGBS chunk. Chunks are then analyzed in parallel in a
24 cluster-computing environment. We calculated the time required by ABBA to analyse chunks of different
25 length (that span from 100 CpGs to 15,000 CpGs) on a single machine with 20 2.3GHz hyper-threaded
26 cores and 32GB of RAM and found that the computational time (seconds) scales with the chunk length
27 (N_{CpG} , number of CpG sites) following the power function: $time \text{ (seconds)} = 0.0045 N_{\text{CpG}}^{1.3985}$ ($R^2 = 0.997$).
28 Depending on the genome length and data dimensionality a complete WGBS analysis ABBA might require
29 days (e.g., it took ~2 weeks to analyse WGBS data in the rat). The total computational time of ABBA
30 analysis can be significantly shortened by splitting the genome into smaller chunks and then assemble the
31 result. The results provided by the “whole-genome” ABBA analysis and “smaller-chunks” ABBA analyses

1 are highly consistent, with no differences in the distribution probabilities obtained with and without
2 splitting the genome in chunks (**Supplementary Figure 4**). Scripts for the pre-processing step are
3 embedded within ABBA at abba.systems-genetics.net

4 **WGBS of rat macrophages.** Bone-marrow derived macrophages (BMDM) were isolated from WKY and
5 LEW rat strains. WGBS libraries were produced as follows: 6µg of genomic DNA was spiked with 10ng of
6 unmethylated cl857 Sam7 lambda DNA (Promega) and sheared using a Covaris System S-series model S2.
7 Sheared DNA was purified and then end-repaired in a 100µl reaction using NEBNext End Repair kit (New
8 England Biolabs) incubated at 20C for 30 minutes. End-repaired DNA was next A-tailed using NEBNext
9 dA-tailing reaction buffer and Klenow Fragment (also New England Biolabs) incubated at 37C for 30
10 minutes and then purified with the MinElute PCR purification kit (Qiagen) in a total final elution volume of
11 28µl. Illumina Early Access Methylation adapter oligos (Illumina) were then ligated to a total of 25µl of the
12 A-tailed DNA sample using NEBNext Quick Ligation Reaction Buffer and Quick T4 DNA ligase (both New
13 England Biolabs) in a reaction volume of 50µl. This mixture was incubated for 30 minutes at 20C prior to
14 gel purification. Bisulphite conversion of 450ng of the purified DNA library was achieved using the Epitect
15 Bisulfite kit (Qiagen) in a total volume of 140µl. Samples were incubated with the following program: 95C
16 for 5 minutes, 60C for 25 minutes, 95C for 5 minutes, 60C for 85 minutes, 95C for 5 minutes, 60C for 175
17 minutes and then 3x repeat of 95C for 5 minutes and 60C for 180 minutes and held at 20C. Treated
18 samples were then purified as per manufacturers instructions. Adapter bound DNA fragments were
19 amplified by a 10-cycle PCR reaction and then purified using Agencourt AMPure XP beads (Beckman
20 Coulter) before gel extraction and quantification using the Agilent Bioanalyzer 2100 Expert High
21 Sensitivity DNA Assay. Then, libraries were quantified using quantitative PCR and then denatured into
22 single stranded fragments. These fragments were then amplified by the Illumina cluster robot and
23 transferred to the HiSeq 2000 for sequencing. WGBS reads were aligned and filtered according to a
24 previously published pipeline (see (Johnson *et al.* 2012) and (Johnson *et al.* 2014)). Briefly, reads were
25 pre-processed by in silico conversion of C bases to T bases in read 1 and G bases to A bases in read 2,
26 followed by clipping of the first base from each read. Pre-processed reads were aligned to the rat
27 reference genome (RGSC3.4) using BWA version 0.6.1 (Li and Durbin 2009) with 3' end quality trimming
28 using a Q score cutoff of 20. Converted and clipped reads 1 and 2 were mapped to two in silico converted
29 versions of the reference sequence, firstly with Cs converted to Ts to allow forward strand mapping, and
30 secondly with Gs converted to As to allow mapping of reverse strand. Aligned reads were filtered by
31 removal of clonal reads, reads with a mapping quality of <20, reads that mapped to both in silico
32 converted forward and reverse strands, and reads with an invalid mapping orientation. We obtained 79.9

1 billion ‘mappable’ bases across both rat strains, with 13.5x (average) coverage in the Lew strain and 17.6x (average) in WKY, where the greatest depth of coverage was observed within CpG islands.

3 Despite ABBA being able to detect methylation changes at all genomic locations we focused only on
4 those methylation changes that occur at CpG sites, and considered CpG sites where at least 4 out of the 8
5 samples contain reads (resulting in a total of 14,976,632 CpG sites genome-wide in BMDM from WKY and
6 LEW rats). DMRs were called with ABBA (see above) using a 5 CpG minimum, a 33% or greater difference
7 in methylation and a 5% FDR threshold. Genomic region annotations and Ensembl gene IDs for the rat
8 reference genome 4 (rn4) were downloaded from the UCSC genome browser. Significant over-
9 representations of genomic features (intron, exons, etc.) were determined empirically from 1,000
10 randomly sampled length and GC-matched regions per DMR. The genes overlapping with DMRs were
11 further annotated and tested for enrichment in Kyoto Encyclopedia of Genes and Genomes (KEGG)
12 pathways using WebGestalt (Wang *et al.* 2013).

13 Identification of enriched transcription factor binding site (TFBS) motifs within the DMRs identified by
14 ABBA was performed using HOMER (Heinz *et al.* 2010). HOMER was used to scan for motifs obtained from
15 the JASPAR 2014 database (Mathelier *et al.* 2014). Threshold used for motifs identification was a p-value
16 of 10^{-4} . Enrichments were calculated by comparing the motifs present in the DMRs against a large set of
17 background sequences ($N = 10^6$) corrected for CpG content.

18 **RNA-seq and ChIP-seq analysis of rat macrophages.** RNA-seq data from BMDM in WKY and LEW
19 strains were retrieved from (Rotival *et al.* 2015) and reanalyzed in the context of WGBS analysis reported
20 here. Briefly, total RNA was extracted from BMDM at day 5 of differentiation in three WKY rats and three
21 LEW rats using Trizol (Invitrogen). 1 μ g of total RNA was used to generate RNA-seq libraries using TruSeq
22 RNA sample preparation kit (Illumina, UK). Libraries were run on a single lane per sample of the HiSeq
23 2000 platform (Illumina) to generate 100bp paired-end reads. An average depth of 72M reads per sample
24 was achieved (minimum 38 M). RNA-seq reads were aligned to the rn4 reference genome using tophat2.
25 The average number of mapped was 67M (minimum 36M) corresponding to an average mapping
26 percentage of 93%. Sequencing and mapping were quality controlled using the FastQC software. Gene-
27 level read counts were computed using HT-Seq-count (Anders *et al.* 2015) with ‘union’ mode and genes
28 with less than 10 aligned reads across all samples were discarded prior to analysis leading to 15,155
29 genes. Differential gene expression analysis between WKY and LEW BMDMs was performed using DESeq2
30 (Love *et al.* 2014) and significantly differentially expressed genes were reported at the 5% FDR level. The
31 visualizations of the expression levels with gene structure were created with DEXSeq (Anders *et al.* 2012).

1 ChIP-seq data from BMDM isolated from the WKY and WKY.LCrgn2 congenic strains (in which the LEW
2 Crgn2 QTL was introgressed onto the WKY background) were retrieved from (Hull *et al.* 2013; Srivastava
3 *et al.* 2013) and re-analyzed with respect to the *Ifitm3* locus. This congenic model (WKY.LCrgn2) has been
4 extensively studied in previous studies where it has been shown that JunD expression levels are
5 significantly higher in WKY when compared with the congenic (Hull *et al.* 2013) and that the canonical
6 binding of AP-1 is significantly greater in WKY compared to WKY.LCrgn2 (Behmoaras *et al.* 2008). Briefly,
7 ChIP was performed with a JunD antibody (Santa Cruz sc74-X) and a negative IgG control (sc-2026). Single
8 read library preparation and high throughput single read sequencing for 36 cycles was carried out on an
9 Illumina Genome Analyser Iix and sequencing of the ChIP-seq libraries was carried out on the high
10 throughput Illumina Genome Analyzer II. Initial data processing was performed using Illumina Real Time
11 Analysis (RTA) v1.6.32 software (equivalent to Illumina Consensus Assessment of Sequence and Variation,
12 CASAVA 1.6) using default settings. Quality filtered reads were then realigned to the rn4 using the
13 Burrows Wheeler Alignment tool v0.5.9 (BWA). Read ends were trimmed if Phred-scaled base quality
14 scores dropped below 20. For the ChIP-seq analysis presented in Figure 3g, differences in JunD binding
15 were assessed only within a 700bp region spanning the *Ifitm3* gene promoter, which included the 600 bp-
16 long DMR identified by ABBA at this locus. ChIP-seq differences were assessed by means of Fisher's exact
17 test on the ChIP-seq counts (normalized for library size) in WKY LCrgn2 and LEW strains, respectively,
18 using a sliding window of 50bp. This locus-specific analysis identified a single 50bp window with
19 differential JunD binding with FET p-value<0.05 that overlapped with JunD TFBS motifs identified by
20 HOMER (see above).

21 **Software and data availability.** ABBA is implemented as a Perl/R program, which is available with
22 instructions for download at abba.systems-genetics.net or via [http://www.mrc-](http://www.mrc-bsu.cam.ac.uk/software/bioinformatics-and-statistical-genomics/)
23 bsu.cam.ac.uk/software/bioinformatics-and-statistical-genomics/. The data is available on Gene
24 Expression Omnibus (GEO), <https://www.ncbi.nlm.nih.gov/geo/>, under the accession number GSE84719.

25 RESULTS

26 We employ a fully Bayesian approach (a Bayesian structured generalized mixed additive model with a
27 latent Gaussian field) which models the random sampling process of the WGBS experiment (the number of
28 methylated/unmethylated reads distributed as non-Gaussian response variable) and where all the
29 unknown quantities are specified by probability distributions. To perform inference ABBA takes
30 advantage of the Integrated Nested Laplace Approximation (INLA) (Rue *et al.* 2009), a new inferential tool
31 for latent Gaussian models. INLA provides approximations to the posterior distribution of the unknowns.

1 These approximations are both very accurate and extremely fast to compute compared to established
2 exact sampling-based methods such as Markov chain Monte Carlo (Gilks *et al.* 1996) (MCMC) or
3 Sequential Monte Carlo (Doucet *et al.* 2001) (SMC). Our new proposed algorithm ABBA is therefore the
4 combination of an approximate inferential procedure with a fully Bayesian model tailored for bisulphite
5 sequencing analysis.

6 ABBA calculates the posterior methylation probability (PMP) at each CpG site based on an estimate of
7 the posterior probability of a smoothed unobserved methylation profile. It also identifies DMRs at a
8 specified FDR by contrasting PMPs across the whole-genome between two groups, e.g. cases and controls.
9 Several intrinsic features of WGBS data are incorporated into ABBA: for instance, the variability in DNA
10 methylation between the (experimental) replicates within each group is modeled through a random effect
11 with a specific within-group variance (**Fig. 1a**). The correlation of DNA methylation patterns is encoded in
12 the latent Gaussian field equation, which reflects the neighborhood structure of the model and
13 automatically adapts to the changes in the underlying data. In particular, the *a priori* correlation between
14 neighbouring CpGs' methylation profiles depends on the distance between them, as it decreases as this
15 distance increases (**Fig. 1b**). Rather than relying on a user-defined value to parameterize it (e.g., kernel
16 bandwidth or window size) or fixing it by an automatic procedure (for instance through an empirical
17 Bayes approach), ABBA assigns a prior distribution on the parameters of the latent Gaussian field
18 equation, thus fully accounting for the uncertainty about these quantities. This specification is key in our
19 model since the data-adaptivity of the degree of smoothing conforms better to the data than assuming
20 fixed values. All these features allow our model to adjust routinely to real-world scenarios, providing an
21 automatic way to describe the WGBS data without requiring any user-defined parameters (Yu and Sun
22 2016b). Full technical details of ABBA algorithm can be found in the Materials and Methods.

23 We benchmarked ABBA and compared it against recently proposed methods (MethylKit (Akalin *et al.*
24 2012), MethylSig (Park *et al.* 2014), DSS/DSS-single (Feng *et al.* 2014; Wu *et al.* 2015), simply DSS
25 hereafter, BSmooth (Hansen *et al.* 2012), metilene (Jühling *et al.* 2015) and the univariate Fisher's exact
26 test (FET)). All methods were run using their default parameterization and for the FET we pooled data
27 from different replicates. To ensure a fair comparison, we used WGBSSuite (Rackham *et al.* 2015) to
28 generate a large number of diverse datasets that were independent of the underlying statistical models of
29 ABBA and of the other methods. Briefly, we simulated *in-silico* datasets to assess the performance of each
30 method under several scenarios, which reflect differences in data integrity and quality of the signal that
31 can occur as a result of biological and experimental phenomenon. The parameters considered were the

1 following: the number of replicates within each group (r), the average read depth per CpG, the level of
2 noise variance (s_0), the methylation probability difference between the two groups ($\Delta meth$) and the
3 switching of methylation status of CpG sites between the two groups (δ) (see **Methods** for details). For
4 each simulated case we generate five replicates and we compared the accuracy of the CpGs called as being
5 contained within DMRs by each technique with the true simulated DMRs. To quantitatively assess the
6 performance of ABBA with respect to competing methods, we evaluated false-positive and false-negative
7 rates of CpG sites and generated receiver operator characteristic (ROC) curves. We focus on the partial
8 area under the ROC curve (or pAUC) at a specificity of 0.75. The pAUC is considered to be more practically
9 relevant than the area under the entire ROC curve (Ma *et al.* 2013) since in typical genomics studies only
10 the features identified at very low false positive rates are selected for further biological validation.

11 All results of the benchmark are detailed in **Supplementary Figures 5-7**. In **Fig 2a** we show
12 representative ROC curves from a specific combination of parameters whilst in **Fig. 2b** we summarize the
13 performance over all combinations of parameters by displaying the best performing method based on its
14 pAUC. Specifically, in **Fig. 2b** the color code in the “benchmark grid” indicates the best performing method
15 for each of the 216 simulated scenarios. For instance, in **Fig. 2a** the top left panel (i) shows the ROC curves
16 for all methods considered under a simulated dataset with $s_0 = 0.1$, $\Delta meth = 30\%$, $r = 1$, average read
17 depth per CpG of 10x and $\delta = 0$. For this combination of parameters we compared the pAUC of each
18 approach, which shows that ABBA is the best performing method. Accordingly, in **Fig. 2b** the square in the
19 grid that represents this parameter set (indicated by (i) in the figure) is coloured black (ABBA). Examples
20 of other ROC curves for specific combinations of parameters are reported in **Fig. 2a** (i-vi) and the
21 corresponding best performing methods are indicated in **Fig. 2b**. In some simulated cases (e.g., with high
22 levels of $\delta = 10\%$) the ROC curves and corresponding pAUC do not distinguish unambiguously the best
23 performing method (e.g., **Fig. 2a** – panel (vi)). In these cases when the pAUC of two methods are very
24 similar ($\pm 1\%$) we report the colours of both methods, e.g., black and red colours in the same square to
25 indicate similar performance of ABBA and DSS (**Fig. 2b**). For the metilene approach (Jühling *et al.* 2015)
26 (which was run using its default parametrization) we noticed that ROC curve analysis was not suitable to
27 compare its performance with other methods. Specifically, for metilene we found that it was not possible to
28 assess both specificity and sensitivity across the wide range of DMRs and scenarios simulated in our
29 study. Representative examples for the ROC curves obtained by running metilene (and other approaches)
30 on the simulated data are provided in **Fig. 2a** and in **Supplementary Figure 8**.

1 Considering all 216 simulated datasets and comparing the pAUCs obtained by each approach across all
2 combinations of parameters, ABBA (black) showed to be the best performing method in 139 (64%) cases
3 (**Fig. 2b-c**). The two other competitive methods were DSS and BSmooth, which show to be the best
4 performing approach only in 26 (12%) and 22 (10%) simulated cases, respectively (**Fig. 2b-c**). In 28
5 (13%) cases different methods showed very similar performance (i.e., pAUCs $\pm 1\%$), and in 17 simulations
6 ABBA and DSS showed to have comparable performance. Looking at the detailed ROC curves reported in
7 **Supplementary Figures 5-7**, we notice that while ABBA was the best method across all simulations (**Fig.**
8 **2c**), its performance diminished for simulated datasets with very small methylation probability difference
9 between the two groups. In particular, for most of the simulated scenarios with $\Delta meth = 20\%$, BSmooth
10 showed very good and robust performance, while DSS was consistently the best performing method when
11 $r = 1$ and $\Delta meth = 20\%$, **Fig. 2b**. However, we highlight that such small differences in DNA methylation
12 (i.e., $\Delta meth \leq 20\%$) are unlikely to have an important biological effect, and the most commonly observed
13 effect sizes for DMR range between 20 and 40%, as previously reported (Ziller *et al.* 2015). In the range
14 $\Delta meth \geq 30\%$, ABBA was the best performing method in 132 (81%) simulations, while DSS was the best
15 performing method only in 10 (6%) simulated cases and, notably, BSmooth was never the best single
16 performing method (BSmooth showed similar performance of ABBA in only one simulated case) (**Fig. 2b**).

17 Specific observations have to be addressed when high levels of errors due to the switching of
18 methylation status of CpG sites between the two groups have been simulated. In these scenarios, it was
19 more difficult to single out a method that outperforms all competing approaches. However, when δ was as
20 high as 10% (i.e., 1 in 10 CpGs is misclassified as unmethylated or vice versa), we observed that ABBA was
21 the best single method in 33 (46%) of 72 simulated scenarios, whereas DSS and BSmooth performed as
22 the best method in 16 (22%) and 7 (10%) of cases, respectively, and in other 10 cases ABBA and DSS have
23 comparable performance. The latter was more apparent when large probability differences between the
24 two groups were simulated ($\Delta meth = 50\%$ or 70%).

25 We then explored whether non-homogeneous, spatially correlated read depth has an effect on ABBA's
26 performance. In order to capture spatially correlated read depth from real data we sampled 5,000
27 contiguous CpGs from WGBS data (generated in rat macrophages, see below and **Methods** for details) and
28 then varied other parameters (r and $\Delta meth$) using WGBSSuite as described above. In these "data-derived"
29 simulated datasets the read depth was correlated with the distance between CpGs (**Supplementary**
30 **Figure 9a**). The results of the benchmark using read depth taken from real data were very similar to those
31 obtained using read depth simulated by means of a Poisson distribution (see **Methods**). Regardless of

1 whether “data-derived” or “Poisson-simulated” read depth was used in our simulations, ABBA was the
2 best performing method to recall DMRs (representative examples are reported in **Supplementary Figure**
3 **9b**). While heterogeneous levels of read depth impact on the single base probability of methylation, the
4 hierarchical model underlying ABBA borrows information across the sequence analyzed, it turns out that
5 ABBA posterior estimates are less sensitive to different levels of the read depth.

6 Taken together our simulation study shows that while individual approaches can be very powerful in
7 detecting DMRs under specific scenarios (notably, DSS with $r = 1$ and BSmooth with $\Delta meth = 20\%$), their
8 performance can vary (and drop) significantly for different choices of the parameters tested in our
9 simulations (at least within the parameter-space considered here). In contrast, we show that, on the
10 whole, ABBA is the best performing method across a large number of parameters’ combination tested and
11 accurately identifies DMRs in the large majority of simulated cases (**Fig. 2c**). Specifically, ABBA’s
12 performance was the highest in the detection of biologically meaningful changes in DNA methylation
13 ($\Delta meth \geq 30\%$) and when little or no errors due to random switching of methylation status of CpG sites
14 between the two groups are present in the data.

15 DNA methylation is emerging as a major contributing factor in several human disorders (Zoghbi and
16 Beaudet 2016), including important autoimmune diseases like systemic lupus erythematosus (SLE) (Wu
17 *et al.* 2016). For instance, differential DNA methylation analysis in CD4+ T cells in lupus patients
18 compared to normal healthy controls identified several genes with known involvement in autoimmunity
19 (Jeffries *et al.* 2011). Here, to illustrate the practical utility of ABBA for differential methylation analysis in
20 disease, we generated WGBS data in an established experimental rat model of crescentic
21 glomerulonephritis (CRGN)(Aitman *et al.* 2006). In this model, we and others have previously shown that
22 susceptibility to CRGN is mediated by macrophages (Behmoaras *et al.* 2008; Page *et al.* 2012); therefore,
23 we assayed CpG methylation at single-nucleotide resolution by WGBS in primary macrophages derived
24 from Wistar Kyoto (WKY) and Lewis (LEW) isogenic rats (two strains discordant for their predisposition
25 to develop CRGN). We used ABBA to carry out genome-wide differential DNA methylation analysis in
26 primary bone-marrow derived macrophages (BMDM) derived from the disease-prone rat strain (WKY, $r =$
27 4) and control strain (LEW, $r = 4$) - see **Methods** for additional details on WGBS data generation and
28 processing. Briefly, in our ABBA analysis of the macrophage methylome, we used the following (default)
29 settings: a minimum of 5 CpG and at least 33% difference in DNA methylation between the disease and
30 control macrophages to identify DMRs. This choice was motivated and supported by data on the local
31 topology of CpG sites in the methylome showing the vast majority of the CpG clusters are in the range of

1 1–11 CpGs (Lövkvist *et al.* 2016) and to increase true positive rate in our DM analysis, following previous
2 assessment and recommendations for methylation analysis using WGBS data (Ziller *et al.* 2015).

3 Using an FDR cutoff of 5%, ABBA identified 1,004 DMRs genome-wide, with 1.07% falling within an
4 annotated CpGI and 6.78% within an annotated CpGS (**Fig. 3a**). For comparative purposes we also used
5 DSS (since this method performed very similarly to ABBA in several simulated cases, **Fig. 2**) to identify
6 DMRs genome-wide, which resulted in only 207 regions with significant differential methylation
7 (uncorrected p-value threshold = 10^{-3} , using the default parameters of DSS). Of the 1,004 DMRs identified
8 by ABBA, 427 overlapped with annotated genes (**Supplementary Table 3**), and there was a significant
9 enrichment for DMRs occurring within 1kb of the gene boundaries (p-value<0.001), within exons (p-
10 value<0.05) and introns (p-value<0.05), **Fig. 3b**. The genes that are within 1kb of a DMR were enriched
11 for pathways relevant to the pathophysiology of CRGN, including MAPK signalling (Ryan *et al.* 2011),
12 Phosphatidylinositol signalling (Wu *et al.* 2014) and Fc gamma R-mediated phagocytosis (Page *et al.*
13 2012) (**Fig. 3c**). For comparison, the 207 DMRs identified by DSS overlapped with 45 genes
14 (**Supplementary Table 4**), which were enriched only for RNA degradation and metabolic pathways. The
15 analysis of real WGBS data by DSS highlighted how the choice of parameters (in this case related to the
16 size of the moving average window in the smoothing procedure) can affect the results. Since the window
17 size in DSS is a user-defined parameter, we performed the analysis with DSS using three different
18 windows (50 bp, 100 bp, 1,000 bp) in addition to the default window size of 500 bp. Each of the four
19 window sizes identified a different number of DMRs, which overlap with different genes (**Supplementary**
20 **Figure 10a**) and have varying distributions of DMR lengths (**Supplementary Figure 10b-e**). The genes
21 identified by DSS when a window of 50 bp is used showed no significant enrichment for pathways, while
22 the results obtained for 100 bp and 1,000 bp windows showed a significant enrichment for RNA
23 degradation. These analyses highlight how the arbitrary choice of parameters related to the degree of
24 smoothing can affect greatly the results of a genome-wide DM analysis as well as the downstream
25 annotation of the genes overlapping with DMRs. In contrast, ABBA automatically adapts to different
26 correlation structures in DNA methylation levels across the genome without requiring any user-defined
27 parameters related to the smoothing procedure.

28 As DNA methylation can affect gene expression by interfering with transcription factor binding, we
29 performed a transcription factor binding site (TFBS) analysis of the DMRs (**Fig. 3d**). This revealed
30 significant enrichment for several TFs, including the ETS transcription factors family and a number of
31 proteins that make the AP-1 TF complex (JUNB, FOS, JUN and JUND), which have been previously linked
32 with CRGN (Behmoaras *et al.* 2008)(Raffetseder *et al.* 2004). To further investigate the potential effect of

1 the changes in DNA methylation identified by ABBA, we carried out differential expression (DE) analysis
2 in macrophages from WKY and LEW rats by RNA-seq (see **Methods** for details). The list of DE genes
3 (n=910, Benjamini–Hochberg (BH)-corrected p-value<0.05) was crosschecked with the genes impacted by
4 DMRs (above), identifying 48 genes with both significant differential methylation and differential
5 expression (**Supplementary Table 5**). We observed the “textbook” model describing DNA methylation
6 regulating transcription via the promoter region (i.e., hypermethylation in the promoter associated with
7 transcriptional repression, see below) as well as widespread methylation changes in the genes body and
8 3’UTR associated with both gene repression and activation. The genes with concordant promoter
9 hypermethylation and transcriptional repression, *Ifitm3*, *Ydjc* and *Cd300Ig* were investigated in more
10 detail since the gene’s promoter is a key regulatory region where the effect of DNA methylation is more
11 clearly understood. We found the biggest change in mRNA expression was in interferon induced
12 transmembrane protein 3 (*Ifitm3*), with mRNA from this gene being almost undetected in unstimulated
13 WKY macrophages (**Fig. 3e**). This observation is consistent with the differential methylation status of the
14 promoter of *Ifitm3*, where the WKY rats had higher levels of methylation than the LEW rats (**Fig. 3f**). To
15 further support the identification of differential methylation at the *Ifitm3* gene we checked whether other
16 methods identified the same DMR. While MethylSig failed to identify significant DMR and BSmooth
17 identified a large and unspecific genomic area as differentially methylated, DSS provides highly consistent
18 results with ABBA, identifying differential methylation at the same region at the *Ifitm3* gene promoter
19 (**Supplementary Figure 11**).

20 We have previously shown that JunD (AP-1) transcription factor is a major determinant of CRGN in
21 WKY rats (Behmoaras *et al.* 2008) and others have shown that AP-1 is methylation sensitive (Ogawa *et al.*
22 2014). Therefore we scanned the DMR (spanning 600 bp) for canonical JunD binding site motifs, and
23 identified three putative regions in the promoter region of *Ifitm3* (**Fig. 3g**). In addition, we re-analyzed
24 ChIP-seq data for JunD transcription factor in BMDM derived from WKY and a congenic strain from LEW
25 (see **Methods** for details). This analysis identified significant differences in JunD binding between WKY
26 and LEW-congenic strain that overlapped with two of the four TFBS identified at the *Ifitm3* promoter (**Fig.**
27 **3g**). The combined evidence provided by our ABBA analysis and RNA-seq/ChIP-seq data therefore
28 suggests that the effect of DNA methylation of the *Ifitm3* gene promoter in WKY rats (prone to develop
29 CRGN) may be restricting the binding of transcription factors such as JunD and, as a consequence, the gene
30 is almost not expressed (<1 TPM) in unstimulated macrophages of WKY rats.

1 DISCUSSION

2 As the cost of genome sequencing technologies continues to drop, it will soon become commonplace to
3 perform comprehensive methylome analyses, using WGBS or other high-throughput techniques that allow
4 the unbiased genome-wide quantification of DNA methylation at a single base-pair resolution. However,
5 high-resolution data generation is only the first step towards the identification of genomic loci and
6 eventually genes with altered methylation levels associated with a given disease, phenotype or
7 developmental stage. The number of DNA methylation datasets available in the public domain is expected
8 to grow; therefore, it becomes necessary to provide the scientific community with analytical tools for a
9 reliable and reproducible identification of differential methylation, and facilitate large epigenome-
10 mapping projects and epigenome-wide association studies (Bock 2012).

11 Beyond statistical power considerations specifically related to the sample size (Rakyan *et al.* 2011) or
12 interpretability of epigenome-wide association studies (Birney *et al.* 2016), our ability to identify
13 accurately changes in DNA methylation localized to specific genomic loci (genes) is also influenced by
14 multiple factors inherently correlated to data quality. These include the within-group heterogeneity, the
15 level of noise, the presence of known genetic covariates (Zhang, 2015) and non-genetic confounding
16 factors (e.g., batch effects) as well as features such as sequencing depth (Ziller *et al.* 2015) or errors due
17 incomplete bisulphite conversion (Genereux *et al.* 2008). Therefore, any analytical tool that can account
18 for all these factors will reduce the number of false positives maximizing the sensitivity and call the
19 regions of interest (i.e., differentially methylated) as accurately as possible. With this in mind, we designed
20 a differential methylation analysis tool (ABBA) that is robust to different experimental and technical
21 variables (see **Fig. 2**), and that adapts automatically to the varying genomic context and local topology of
22 CpG sites affecting methylation levels. In particular, the automatic adaptation to different correlation
23 structures in CpG methylation levels (without requiring user-defined parameters about the degree of
24 smoothing) as well as the ability of modelling its decay as the function of the genomic distances between
25 CpGs allow ABBA to adapt routinely to methylation changes that occurs with different scales and non-
26 uniform rates across the genome. The importance of the genomic context in the methylome and the local
27 topology of CpG sites have been recently investigated, showing, amongst other features, that methylation
28 at small CpG clusters is more likely to induce stable changes in DNA methylation (Lövkvist *et al.* 2016).

29 From a user's perspective, ABBA treats WGBS-seq data in a general way with no specification of
30 parameters related to the level of data smoothing (such as window size or kernel bandwidth), thus
31 allowing for a great deal of automation. This also facilitates the WGBS analysis when the values of the

1 parameter settings (that may largely affect the accuracy of DM identification) are not known. Our fully
2 Bayesian approach can be also easily modified to include covariates and non-genetic confounding factors
3 through random effects, beyond the replicates level. It also allows the specification of covariates that are
4 informative about the methylation profiles by adding prior biological information to the linear predictor
5 μ_{ig} in (6). While these alterations can be done in our model with a simple modification of the code and
6 with negligible further computational costs, non-parametric smoothing techniques (spline- (Hansen *et al.*
7 2012), kernel (Hebestreit *et al.* 2013)- and moving average-based smoothing (Feng *et al.* 2014)) do not
8 possess the same straightforward flexibility nor alternative approaches based on Hidden Markov Models
9 (Yu and Sun 2016a), (Kuan and Chiang 2012), (Sun and Yu 2016).

10 Our extensive simulation studies (**Fig. 2**) and differential DNA methylation analysis in
11 glomerulonephritis (**Fig. 3**) showed that ABBA is a powerful approach for the identification of DMRs from
12 WGBS single-base pair resolution methylation data. While individual methods such as BSmooth (Hansen
13 *et al.* 2012) or DSS (Feng *et al.* 2014; Wu *et al.* 2015) showed a very good power to detect DMRs under
14 specific scenarios and conditions, ABBA retained a high degree of robustness of the results with respect to
15 a wider range of factors (parameters) affecting WGBS data integrity and quality, including sequencing
16 coverage, number of replicates or different noise structures. This is particularly appealing in cases when
17 considerable efforts have been expended toward generation of large-scale WGBS data from heterogeneous
18 systems, e.g., the ENCODE project (Bernstein *et al.* 2012), and data quality can vary across experimental
19 conditions and laboratories. As proof of concept of ABBA's application to real data analysis, we used an
20 established experimental model system of glomerulonephritis (Aitman *et al.* 2006) to identify changes in
21 DNA methylation associated with disease. In this, we employed ABBA to analyze ~15 million CpG sites
22 genome-wide in primary bone-marrow derived macrophages derived from WKY and LEW rats and
23 identified >1,000 significant DMRs at 5% FDR level. A comparative analysis using DSS (the most
24 competitive approach from our simulation study) did not provide the same level of biological insight both
25 in terms of significant pathway enrichments and in robustly identifying DMRs across user-defined
26 parameters. To highlight this point, we showed how the results of DSS were greatly affected by the choice
27 of the window size.

28 Furthermore, we have shown how integrating the DMR results provided by ABBA with other 'omics'
29 data (RNA-seq and ChIP-seq generated in the same experimental system), enabled us to generate new
30 hypotheses for the mechanism underpinning the disease, revealing a candidate gene (*Ifitm3*) for the
31 susceptibility to glomerulonephritis. These findings on *Ifitm3* in rat glomerulonephritis merit further
32 discussion. *Ifitm3* has a known role in viral resistance, a central part of innate immunity, and is inducible

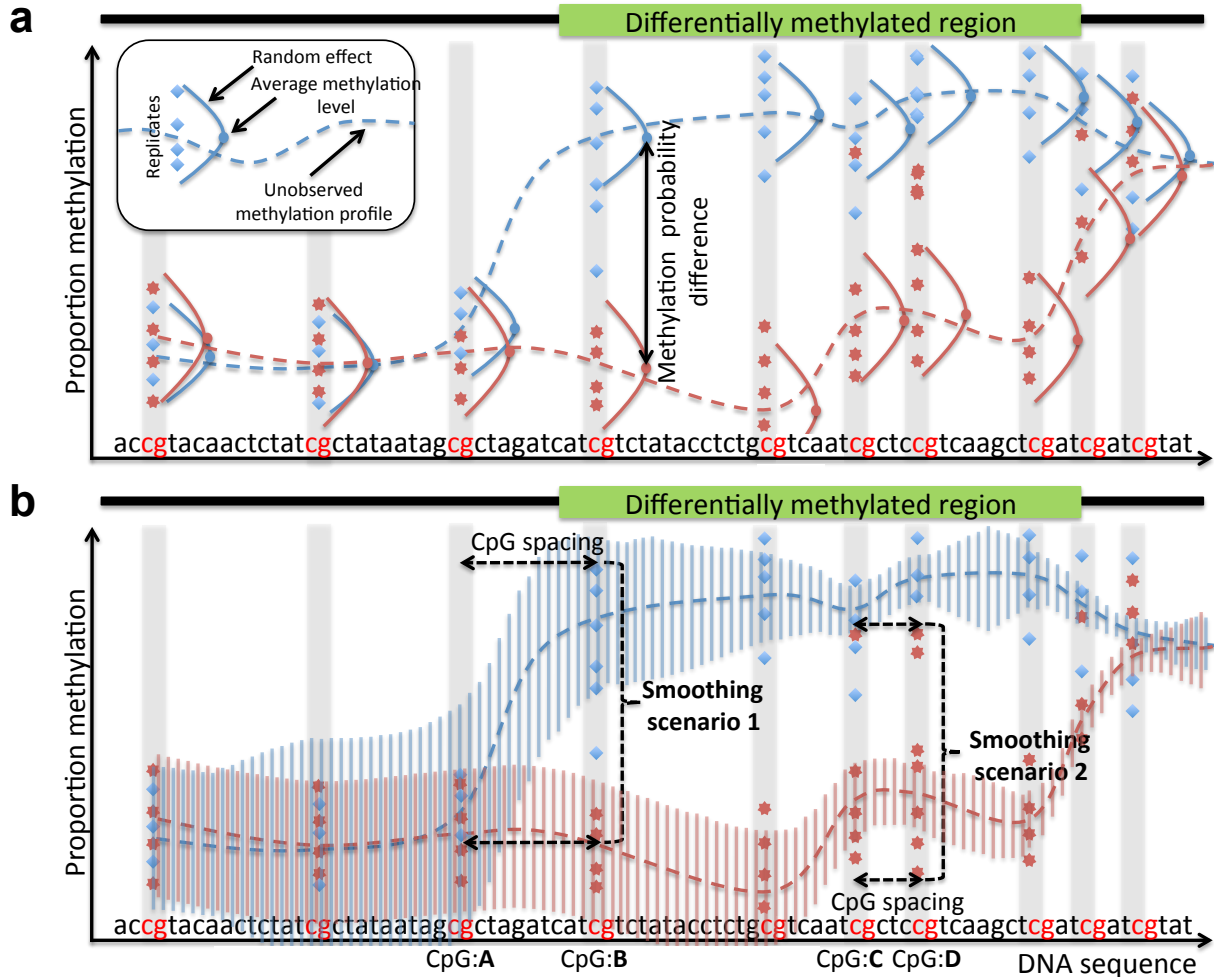
1 by both interferon (IFN) types I and II (Everitt *et al.* 2012). Notably, type II IFN signaling has been
2 implicated in the pathogenesis of nephrotoxic nephritis and other “planted” antigen models of CRGN
3 (Kitching *et al.* 2004), although DNA methylation has not previously been examined in this context. With
4 regards to type I IFN, recent genome-wide DNA methylation analysis of T-cells, B-cells and monocytes has
5 shown that patients with SLE, a frequent autoimmune cause of CRGN, have severe hypomethylation near
6 to genes involved in type I IFN signaling (Absher *et al.* 2013). In addition, DNA methylation alterations in
7 IFN-related genes, including *Ifitm3*, have been previously observed and proposed to contribute to the
8 pathogenesis of other autoimmune diseases such as primary Sjögren's syndrome (Gottenberg *et al.* 2006).
9 Regarding the role of *Ifitm3* gene, it has been shown to directly interact *in vivo* and *in vitro*, with
10 osteopontin, a matricellular protein, whose transcription is mediated by the AP-1 TF family (El-Tanani *et*
11 *al.* 2010). Furthermore, osteopontin has been also previously associated with SLE (Rullo *et al.* 2013) and
12 ANCA-associated vasculitis (Lorenzen *et al.* 2010) another frequent cause of CRGN. Therefore, our ABBA
13 analysis of WGBS data in primary macrophages from a rat model of CRGN allowed us to propose an AP-1-
14 mediated role for *Ifitm3* in glomerulonephritis. While a role for IFN-signaling genes in autoimmune
15 disease has been previously suggested, our findings on methylation alteration of the *Ifitm3* gene
16 associated with glomerulonephritis in the rat might suggest future directions for the study of the
17 pathogenesis and to develop treatments of CRGN.

18 In a wider context, the role of methylation is dependent on the location with respect to the gene body
19 and regulation functions. Methylation in a CpGI-depleted promoter, such as the promoter region of *Ifitm3*
20 gene (according to UCSC genome browser (RN4)), is associated with repression that maybe due to
21 interference with transcription factor binding. Conversely, methylation in the gene body is positively
22 associated with active transcription as methylation can be caused by transcriptional elongation
23 (Schübeler 2015). Methylation within a gene body can also act as an insulator for repetitive and
24 transposable elements or distal intronic enhancers, on which the methylation would have no regulatory
25 effect on the gene in which it resides (Jones 2012). Given the complexity of these regulatory functions of
26 methylation, the ability of our approach to accurately identify changes in DNA methylation that are
27 localized to specific regions is likely to facilitate our understanding of the complex relationships between
28 methylation and gene regulation. As exemplified by our integrative analysis of the of the *Ifitm3* locus, we
29 anticipate that the ABBA results for differential DNA methylation should be integrated with additional
30 transcriptional and epigenetic data in order to better define hypotheses on specific regulatory
31 mechanisms.

1 In summary, we show how ABBA provides a flexible and user-friendly automatic framework for the
2 identification of differential methylation that is robust across a wide range of experimental parameters, an
3 approach that we have also applied to identify changes in macrophage DNA methylation in
4 glomerulonephritis.

5

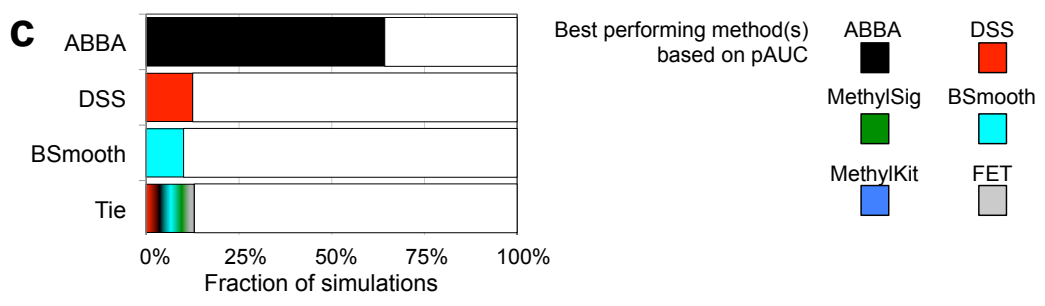
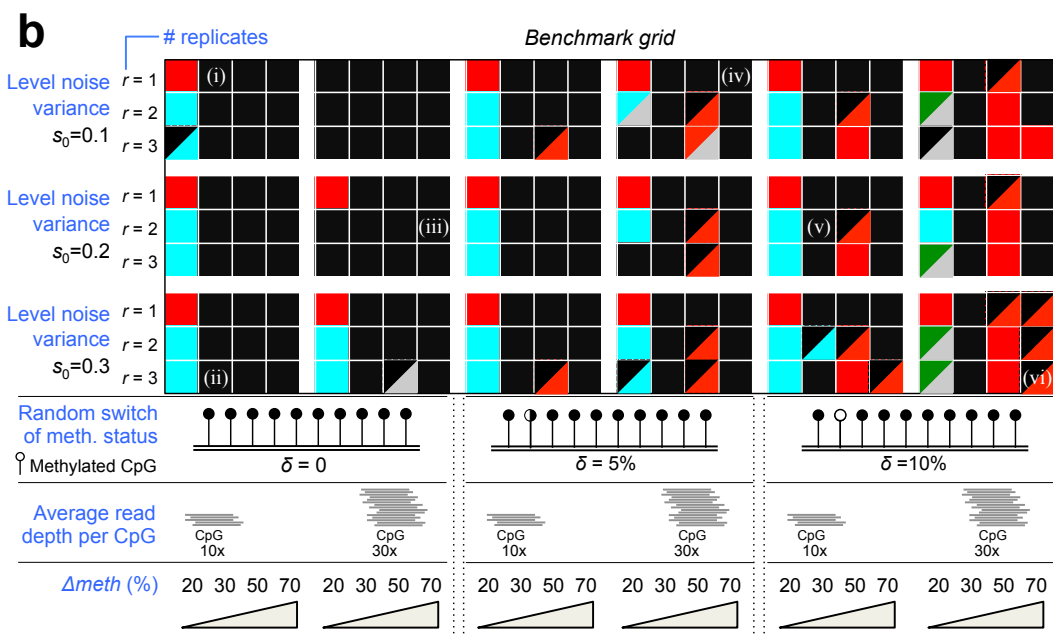
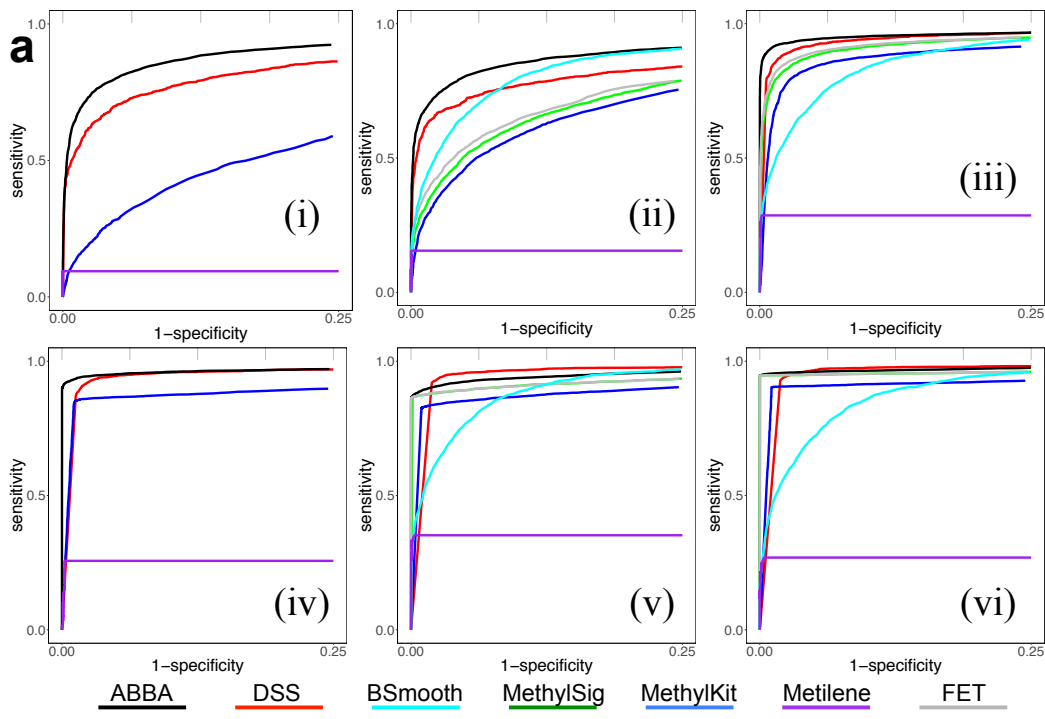
1 FIGURES



2

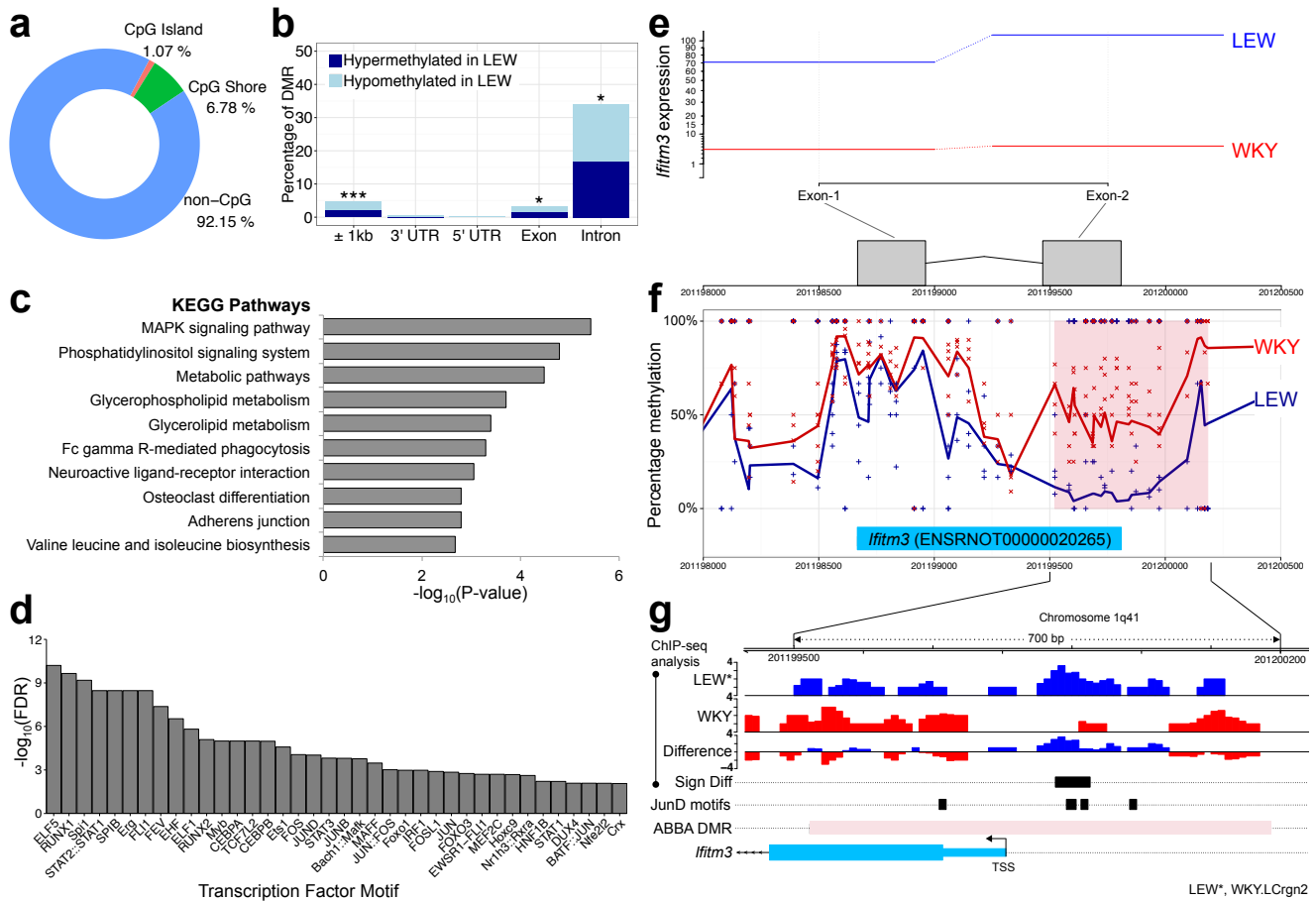
3 **Figure 1. ABBA model.** ABBA estimates the unobserved methylation profiles, i.e. the DNA average
 4 methylation levels across replicates, of two groups from WGBS data (blue diamonds and read stars). (a) A
 5 random effect accounts for the variability of experimental replicates. At each CpG the methylation
 6 probability difference is the difference between the methylation profile of the two groups (blue and red
 7 dots). (b) The methylation profiles of each group are smoothed by a latent Gaussian field that
 8 probabilistically connects them (dotted lines). In particular “Smoothing scenario 1” shows that if a large
 9 spacing (distance) between two consecutive CpGs (CpG:A and CpG:B) exists, the methylation profile at
 10 CpG:B does not depend on the previous one at CpG:A (blue dotted line). The opposite happens in
 11 “Smoothing scenario 2” where the methylation profile at CpG:D is largely influenced by the previous one
 12 at CpG:C (red dotted line) despite some high levels of methylation (red stars) which are treated by ABBA
 13 as outliers. The degree of the smoothing, i.e. the correlation between DNA methylation profiles, is

1 controlled automatically by the marginal variance of the Latent Gaussian Field (blue and red vertical
2 bars): the correlation is higher (lower) when the variance is small (large). On the other hand, the variance
3 decreases as the distance between neighbouring CpGs' decreases (Smoothing scenario 2) while increases
4 as the distance increases (Smoothing scenario 1).



1

1 **Figure 2. Benchmarking results.** (a) ROC curves for selected combinations of parameters: (i) $s_0 = 0.1$,
2 $\Delta meth = 30\%$, $r = 1$, average read depth per CpG of 10x, $\delta = 0$; (ii) $s_0 = 0.3$, $\Delta meth = 30\%$, $r = 3$, average
3 read depth per CpG of 10x, $\delta = 0$; (iii) $s_0 = 0.2$, $\Delta meth = 70\%$, $r = 2$, average read depth per CpG of 30x, $\delta =$
4 0; (iv) $s_0 = 0.1$, $\Delta meth = 70\%$, $r = 1$, average read depth per CpG of 30x, $\delta = 5\%$; (v) $s_0 = 0.2$, $\Delta meth = 30\%$, r
5 =2, average read depth per CpG of 10x, $\delta = 10\%$; (vi) $s_0 = 0.3$, $\Delta meth = 70\%$, $r = 3$, average read depth per
6 CpG of 30x, $\delta = 10\%$. For each of this combination of parameters, the corresponding best method based
7 on its pAUC is indicated in the benchmark grid below. In (i) and (iv) ROC curves are reported only for the
8 methods that can analyze WGBS data generated from one biological sample. (b) Global snapshot of the
9 method's performance across 216 simulated datasets. A given combination of parameters is indicated by a
10 square in the benchmark grid, and for each square we calculated the pAUC for each method and
11 determined which method had the overall best pAUC (i.e., $pAUC_{method_1} > pAUC_{method_2}$). Colours in the
12 benchmark grid indicate which method had the best performance. When pAUC of two methods are similar
13 ($\pm 1\%$) we report the colours of both methods (e.g., black and red colours in the same square indicate
14 similar performance of ABBA and DSS). The six selected combination of parameters for which the ROC
15 curves are reported in panel (a) are indicated within the benchmark grid: (i, ii, iii, iv, v and vi). All ROC
16 curves are reported in **Supplementary Figures 5-7**. (c) For the three best performing methods (ABBA,
17 DSS and BSmooth) we report the percentage of simulated scenarios in which each method resulted to be
18 the best based on the pAUC comparison. "Tie" indicates the proportion of simulated scenarios in which the
19 pAUCs of any two methods were similar (i.e., pAUCs $\pm 1\%$) and it was not possible to single out a single
20 best performing approach.



1

2 **Figure 3. ABBA analysis of WGBS in rat macrophages.** (a) CpG-based annotation 1,004 DMR between
 3 WKY and LEW macrophages showing significantly higher proportions of CpGI and CpGS than those that
 4 would be expected by chance (p-value<0.009 for CpGI and p-value<0.001 for CpGS, respectively, obtained
 5 by 1,000 randomly sampled datasets of 1,004 CpG-matched regions). (b) Proportions of DMRs in different
 6 genomic features of overlapping genes. Feature annotation was retrieved from UCSC genome browser
 7 (RN4). (c) KEGG pathway enrichment for the genes overlapping with DMRs. Only significant pathways are
 8 reported (FDR<1%). (d) Enrichment for the TFBS within the DMRs was when compared to CG matched
 9 regions of the genome (FDR<0.05). (e) RNA-seq analysis in WKY and LEW macrophages shows lack of
 10 *Ifitm3* expression in WKY rats. (f) Percentage methylation at each CpG in WKY (crosses) and LEW (plus)
 11 and smoothed average methylation profiles by ABBA. The pink box highlights the significant DMR
 12 identified by ABBA (FDR<5%). (g) ChIP-seq analysis for JunD in LEW.LCrgn2 (LEW*) and WKY
 13 macrophages identified a single region with differential binding of JunD (p-value<0.05, Sign Diff
 14 row, black box). Units on the y-axis refer to relative ChIP-seq coverage with respect to the control. This region

1 overlapped with two (out of four) JunD binding sites motifs identified within the gene promoter (± 500 bp
2 around the TSS). ABBA DMR, differentially methylated region identified by ABBA. TSS, transcription start
3 site. *, p-value<0.05, ***, p-value<0.001

4

1 **ACKNOWLEDGEMENTS**

2 The authors are thankful to the two anonymous referees whose meticulous attention to their refereeing
3 task has resulted in substantial improvements in our presentation.

4

5 **FUNDING**

6 This research was funded by Engineering and Physical Sciences Research Council Grant EP/K030760/1
7 (L.B.), The Alan Turing Institute under the EPSRC grant EP/N510129/1 (L.B., P.D.), Royal Society
8 IE110977 (L.B., P.D.), European Union (European Social Fund - ESF), Greek national funds through the
9 Operational Program "Education and Lifelong Learning" of the National Strategic Reference Framework
10 (NSRF), project ARISTEIA (P.D.), Duke-NUS Medical School and Singapore Ministry of Health (O.J.L.R.,
11 E.P.), a Medical Research Council Chain-Florey fellowship (T.O.), the Medical Research Council
12 (MR/M004716/1 to J.B. and E.P.) and by Kidney Research UK - RP9/2013 (J.B.). The funders had no role in
13 study design, data collection and analysis, decision to publish, or preparation of the manuscript.

14

15 **AUTHOR CONTRIBUTIONS**

16 L.B. and E.P. initiated, directed and supervised the project. P.D. and L.B. conceived the statistical model
17 and the computational approach. P.D., E.V. and L.B. wrote the initial algorithm that was further developed
18 by O.J.L.R. and L.B. to the presented approach. T.O. and E.P. generated WGBS data in the rat. S.R.L., O.J.L.R.
19 and E.P. carried out analysis of WGBS and RNA-seq data in the rat and interpreted the results. N.H. and
20 P.K.S. carried out ChIP-seq and TFBS analyses. J.B. provided RNA-seq and ChIP-seq data in the rat. O.J.L.R.,
21 L.B. and E.P. wrote the manuscript with input from all authors. All of the authors read and approved the
22 final manuscript.

23

1 REFERENCES

- 2 Absher D. M., Li X., Waite L. L., Gibson A., Roberts K., Edberg J., Chatham W. W., Kimberly R. P., 2013
3 Genome-wide DNA methylation analysis of systemic lupus erythematosus reveals persistent
4 hypomethylation of interferon genes and compositional changes to CD4+ T-cell populations. *PLoS*
5 *Genet.* **9**: e1003678.
- 6 Äijö T., Huang Y., Mannerström H., Chavez L., Tsagaratou A., Rao A., Lähdesmäki H., 2016 A probabilistic
7 generative model for quantification of DNA modifications enables analysis of demethylation
8 pathways. *Genome Biol.* **17**: 49.
- 9 Aitman T. J., Dong R., Vyse T. J., Norsworthy P. J., Johnson M. D., Smith J., Mangion J., Robertson-Lowe C.,
10 Marshall A. J., Petretto E., Hodges M. D., Bhangal G., Patel S. G., Sheehan-Rooney K., Duda M., Cook P.
11 R., Evans D. J., Domin J., Flint J., Boyle J. J., Pusey C. D., Cook H. T., 2006 Copy number polymorphism in
12 *Fcgr3* predisposes to glomerulonephritis in rats and humans. *Nature* **439**: 851–5.
- 13 Akalin A., Kormaksson M., Li S., Garrett-Bakelman F. E., Figueroa M. E., Melnick A., Mason C. E., 2012
14 methylKit: a comprehensive R package for the analysis of genome-wide DNA methylation profiles.
15 *Genome Biol.* **13**: R87.
- 16 Anders S., Reyes A., Huber W., 2012 Detecting differential usage of exons from RNA-seq data. *Genome Res.*
17 **22**: 2008–17.
- 18 Anders S., Pyl P. T., Huber W., 2015 HTSeq—a Python framework to work with high-throughput
19 sequencing data. *Bioinformatics* **31**: 166–169.
- 20 Behmoaras J., Bhangal G., Smith J., McDonald K., Mutch B., Lai P. C., Domin J., Game L., Salama A., Foxwell B.
21 M., Pusey C. D., Cook H. T., Aitman T. J., 2008 *Jund* is a determinant of macrophage activation and is
22 associated with glomerulonephritis susceptibility. *Nat. Genet.* **40**: 553–9.
- 23 Bell J. T., Pai A. A., Pickrell J. K., Gaffney D. J., Pique-Regi R., Degner J. F., Gilad Y., Pritchard J. K., 2011 DNA
24 methylation patterns associate with genetic and gene expression variation in HapMap cell lines.
25 *Genome Biol.* **12**: R10.
- 26 Bernstein B. E., Birney E., Dunham I., Green E. D., Gunter C., Snyder M., 2012 An integrated encyclopedia of
27 DNA elements in the human genome. *Nature* **489**: 57–74.
- 28 Birney E., Smith G. D., Greally J. M., 2016 Epigenome-wide Association Studies and the Interpretation of
29 Disease -Omics. *PLoS Genet.* **12**: e1006105.
- 30 Bock C., 2012 Analysing and interpreting DNA methylation data. *Nat. Rev. Genet.* **13**: 705–19.
- 31 Broët P., Lewin A., Richardson S., Dalmaso C., Magdelenat H., 2004 A mixture model-based strategy for
32 selecting sets of genes in multiclass response microarray experiments. *Bioinformatics* **20**: 2562–71.
- 33 Chen Z. -x., Riggs A. D., 2011 DNA Methylation and Demethylation in Mammals. *J. Biol. Chem.* **286**: 18347–
34 18353.
- 35 Deaton A. M., Bird A., 2011 CpG islands and the regulation of transcription. *Genes Dev.* **25**: 1010–22.
- 36 Dempster A. P. A., Laird N. M. N., Rubin D. D. B., 1977 Maximum likelihood from incomplete data via the
37 EM algorithm. *J. R. Stat. Soc. Ser. B* **39**: 1–38.
- 38 Doucet a, Freitas N. De, Gordon N., 2001 *Sequential Monte Carlo Methods in Practice*.
- 39 Efron B., 2008 Simultaneous inference: When should hypothesis testing problems be combined? *Ann.*
40 *Appl. Stat.* **2**: 197–223.

- 1 El-Tanani M. K., Jin D., Campbell F. C., Johnston P. G., 2010 Interferon-induced transmembrane 3 binds
2 osteopontin in vitro: expressed in vivo IFITM3 reduced OPN expression. *Oncogene* **29**: 752–762.
- 3 Everitt A. R., Clare S., Pertel T., John S. P., Wash R. S., Smith S. E., Chin C. R., Feeley E. M., Sims J. S., Adams D.
4 J., Wise H. M., Kane L., Goulding D., Digard P., Anttila V., Baillie J. K., Walsh T. S., Hume D. A., Palotie A.,
5 Xue Y., Colonna V., Tyler-Smith C., Dunning J., Gordon S. B., GenISIS Investigators, MOSAIC
6 Investigators, Smyth R. L., Openshaw P. J., Dougan G., Brass A. L., Kellam P., 2012 IFITM3 restricts the
7 morbidity and mortality associated with influenza. *Nature* **484**: 519–23.
- 8 Feng H., Conneely K. N., Wu H., 2014 A Bayesian hierarchical model to detect differentially methylated loci
9 from single nucleotide resolution sequencing data. *Nucleic Acids Res.* **42**: e69.
- 10 Genereux D. P., Johnson W. C., Burden A. F., Stöger R., Laird C. D., 2008 Errors in the bisulfite conversion of
11 DNA: modulating inappropriate- and failed-conversion frequencies. *Nucleic Acids Res.* **36**: e150.
- 12 Gilks W. R., Richardson S., Spiegelhalter D. J., 1996 *Markov Chain Monte Carlo in Practice*.
- 13 Gottenberg J.-E., Cagnard N., Lucchesi C., Letourneur F., Mistou S., Lazure T., Jacques S., Ba N., Ittah M.,
14 Lepajolec C., Labetoulle M., Ardizzone M., Sibilia J., Fournier C., Chiocchia G., Mariette X., 2006
15 Activation of IFN pathways and plasmacytoid dendritic cell recruitment in target organs of primary
16 Sjögren's syndrome. *Proc. Natl. Acad. Sci. U. S. A.* **103**: 2770–5.
- 17 Gutierrez-Arcelus M., Lappalainen T., Montgomery S. B., Buil A., Ongen H., Yurovsky A., Bryois J., Giger T.,
18 Romano L., Planchon A., Falconnet E., Bielser D., Gagnebin M., Padioleau I., Borel C., Letourneau A.,
19 Makrythanasis P., Guipponi M., Gehrig C., Antonarakis S. E., Dermitzakis E. T., 2013 Passive and active
20 DNA methylation and the interplay with genetic variation in gene regulation. *Elife* **2**: e00523.
- 21 Gutierrez-Arcelus M., Ongen H., Lappalainen T., Montgomery S. B., Buil A., Yurovsky A., Bryois J., Padioleau
22 I., Romano L., Planchon A., Falconnet E., Bielser D., Gagnebin M., Giger T., Borel C., Letourneau A.,
23 Makrythanasis P., Guipponi M., Gehrig C., Antonarakis S. E., Dermitzakis E. T., 2015 Tissue-specific
24 effects of genetic and epigenetic variation on gene regulation and splicing. *PLoS Genet.* **11**: e1004958.
- 25 Hansen K. D., Langmead B., Irizarry R. A., 2012 BSmooth: from whole genome bisulfite sequencing reads to
26 differentially methylated regions. *Genome Biol.* **13**: R83.
- 27 Harris R. A., Wang T., Coarfa C., Nagarajan R. P., Hong C., Downey S. L., Johnson B. E., Fouse S. D., Delaney A.,
28 Zhao Y., Olshen A., Ballinger T., Zhou X., Forsberg K. J., Gu J., Echipare L., O'Geen H., Lister R., Pelizzola
29 M., Xi Y., Epstein C. B., Bernstein B. E., Hawkins R. D., Ren B., Chung W.-Y., Gu H., Bock C., Gnirke A.,
30 Zhang M. Q., Haussler D., Ecker J. R., Li W., Farnham P. J., Waterland R. A., Meissner A., Marra M. A.,
31 Hirst M., Milosavljevic A., Costello J. F., 2010 Comparison of sequencing-based methods to profile
32 DNA methylation and identification of monoallelic epigenetic modifications. *Nat. Biotechnol.* **28**:
33 1097–105.
- 34 Hebestreit K., Dugas M., Klein H.-U., 2013 Detection of significantly differentially methylated regions in
35 targeted bisulfite sequencing data. *Bioinformatics* **29**: 1647–53.
- 36 Heinz S., Benner C., Spann N., Bertolino E., Lin Y. C., Laslo P., Cheng J. X., Murre C., Singh H., Glass C. K., 2010
37 Simple combinations of lineage-determining transcription factors prime cis-regulatory elements
38 required for macrophage and B cell identities. *Mol. Cell* **38**: 576–89.
- 39 Hull R. P., Srivastava P. K., D'Souza Z., Atanur S. S., Mechta-Grigoriou F., Game L., Petretto E., Cook H. T.,
40 Aitman T. J., Behmoaras J., 2013 Combined CHIP-Seq and transcriptome analysis identifies AP-1/JunD
41 as a primary regulator of oxidative stress and IL-1 β synthesis in macrophages. *BMC Genomics* **14**: 92.
- 42 Jeffries M. A., Dozmorov M., Tang Y., Merrill J. T., Wren J. D., Sawalha A. H., 2011 Genome-wide DNA

- 1 methylation patterns in CD4+ T cells from patients with systemic lupus erythematosus. *Epigenetics*
2 **6**: 593–601.
- 3 Johnson M. D., Mueller M., Game L., Aitman T. J., 2012 Single Nucleotide Analysis of Cytosine Methylation
4 by Whole-Genome Shotgun Bisulfite Sequencing. In: *Current Protocols in Molecular Biology*, John
5 Wiley & Sons, Inc., Hoboken, NJ, USA, p. Unit21.23.
- 6 Johnson M. D., Mueller M., Adamowicz-Brice M., Collins M. J., Gellert P., Maratou K., Srivastava P. K., Rotival
7 M., Butt S., Game L., Atanur S. S., Silver N., Norsworthy P. J., Langley S. R., Petretto E., Pravenec M.,
8 Aitman T. J., 2014 Genetic analysis of the cardiac methylome at single nucleotide resolution in a
9 model of human cardiovascular disease. (R Mott, Ed.). *PLoS Genet.* **10**: e1004813.
- 10 Jones P. A., 2012 Functions of DNA methylation: islands, start sites, gene bodies and beyond. *Nat. Rev.*
11 *Genet.* **13**: 484–92.
- 12 Jühling F., Kretzmer H., Bernhart S. H., Otto C., Stadler P. F., Hoffmann S., 2015 metilene: Fast and sensitive
13 calling of differentially methylated regions from bisulfite sequencing data. *Genome Res.:*
14 *gr.196394.115-*.
- 15 Kitching A. R., Turner A. L., Semple T., Li M., Edgton K. L., Wilson G. R., Timoshanko J. R., Hudson B. G.,
16 Holdsworth S. R., 2004 Experimental autoimmune anti-glomerular basement membrane
17 glomerulonephritis: a protective role for IFN-gamma. *J. Am. Soc. Nephrol.* **15**: 1764–74.
- 18 Kuan P. F., Chiang D. Y., 2012 Integrating prior knowledge in multiple testing under dependence with
19 applications to detecting differential DNA methylation. *Biometrics* **68**: 774–83.
- 20 Lea A. J., Tung J., Zhou X., 2015 A Flexible, Efficient Binomial Mixed Model for Identifying Differential DNA
21 Methylation in Bisulfite Sequencing Data (D Absher, Ed.). *PLOS Genet.* **11**: e1005650.
- 22 Li H., Durbin R., 2009 Fast and accurate short read alignment with Burrows-Wheeler transform.
23 *Bioinformatics* **25**: 1754–1760.
- 24 Lorenzen J., Lovric S., Krämer R., Haller H., Haubitz M., 2010 Osteopontin in antineutrophil cytoplasmic
25 autoantibody-associated vasculitis: relation to disease activity, organ manifestation and
26 immunosuppressive therapy. *Ann. Rheum. Dis.* **69**: 1169–71.
- 27 Love M. I., Huber W., Anders S., 2014 Moderated estimation of fold change and dispersion for RNA-seq
28 data with DESeq2. *Genome Biol.* **15**: 550.
- 29 Lövkvist C., Dodd I. B., Sneppen K., Haerter J. O., 2016 DNA methylation in human epigenomes depends on
30 local topology of CpG sites. *Nucleic Acids Res.* **44**: 5123–32.
- 31 Ma H., Bandos A. I., Rockette H. E., Gur D., 2013 On use of partial area under the ROC curve for evaluation
32 of diagnostic performance. *Stat. Med.* **32**: 3449–58.
- 33 Mathelier A., Zhao X., Zhang A. W., Parcy F., Worsley-Hunt R., Arenillas D. J., Buchman S., Chen C., Chou A.,
34 Ienasescu H., Lim J., Shyr C., Tan G., Zhou M., Lenhard B., Sandelin A., Wasserman W. W., 2014 JASPAR
35 2014: an extensively expanded and updated open-access database of transcription factor binding
36 profiles. *Nucleic Acids Res.* **42**: D142–7.
- 37 Ogawa C., Tone Y., Tsuda M., Peter C., Waldmann H., Tone M., 2014 TGF- β -Mediated Foxp3 Gene Expression
38 Is Cooperatively Regulated by Stat5, Creb, and AP-1 through CNS2. *J. Immunol.* **192**: 475–483.
- 39 Page T. H., D'Souza Z., Nakanishi S., Serikawa T., Pusey C. D., Aitman T. J., Cook H. T., Behmoaras J., 2012
40 Role of novel rat-specific Fc receptor in macrophage activation associated with crescentic
41 glomerulonephritis. *J. Biol. Chem.* **287**: 5710–9.

- 1 Park Y., Figueroa M. E., Rozek L. S., Sartor M. A., 2014 MethylSig: a whole genome DNA methylation
2 analysis pipeline. *Bioinformatics* **30**: 2414–22.
- 3 Rackham O. J. L., Dellaportas P., Petretto E., Bottolo L., 2015 WGBSSuite: simulating whole-genome
4 bisulphite sequencing data and benchmarking differential DNA methylation analysis tools.
5 *Bioinformatics* **31**: 2371–3.
- 6 Raffetseder U., Wernert N., Ostendorf T., Roeyen C. van, Rauen T., Behrens P., Floege J., Mertens P. R., 2004
7 Mesangial cell expression of proto-oncogene Ets-1 during progression of mesangioproliferative
8 glomerulonephritis. *Kidney Int.* **66**: 622–32.
- 9 Rakyan V. K., Down T. A., Balding D. J., Beck S., 2011 Epigenome-wide association studies for common
10 human diseases. *Nat. Rev. Genet.* **12**: 529–41.
- 11 Robinson M. D., Kahraman A., Law C. W., Lindsay H., Nowicka M., Weber L. M., Zhou X., 2014 Statistical
12 methods for detecting differentially methylated loci and regions. *Front. Genet.* **5**: 324.
- 13 Rotival M., Ko J.-H., Srivastava P. K., Kerloc’h A., Montoya A., Mauro C., Faull P., Cutillas P. R., Petretto E.,
14 Behmoaras J., 2015 Integrating phosphoproteome and transcriptome reveals new determinants of
15 macrophage multinucleation. *Mol. Cell. Proteomics* **14**: 484–98.
- 16 Rue Hå., Martino S., 2007 Approximate Bayesian inference for hierarchical Gaussian Markov random field
17 models. *J. Stat. Plan. Inference* **137**: 3177–3192.
- 18 Rue H., Martino S., Chopin N., 2009 Approximate Bayesian inference for latent Gaussian models by using
19 integrated nested Laplace approximations. *J. R. Stat. Soc. Ser. B* **71**: 319–392.
- 20 Rullo O. J., Woo J. M. P., Parsa M. F., Hoftman A. D. C., Maranian P., Elashoff D. A., Niewold T. B., Grossman J.
21 M., Hahn B. H., McMahan M., McCurdy D. K., Tsao B. P., 2013 Plasma levels of osteopontin identify
22 patients at risk for organ damage in systemic lupus erythematosus. *Arthritis Res. Ther.* **15**: R18.
- 23 Ryan J., Ma F. Y., Kanellis J., Delgado M., Blease K., Nikolic-Paterson D. J., 2011 Spleen tyrosine kinase
24 promotes acute neutrophil-mediated glomerular injury via activation of JNK and p38 MAPK in rat
25 nephrotoxic serum nephritis. *Lab. Invest.* **91**: 1727–38.
- 26 Schübeler D., 2015 Function and information content of DNA methylation. *Nature* **517**: 321–326.
- 27 Srivastava P. K., Hull R. P., Behmoaras J., Petretto E., Aitman T. J., 2013 JunD/AP1 regulatory network
28 analysis during macrophage activation in a rat model of crescentic glomerulonephritis. *BMC Syst.*
29 *Biol.* **7**: 93.
- 30 Sun W., Tony Cai T., 2009 Large-scale multiple testing under dependence. *J. R. Stat. Soc. Ser. B (Statistical*
31 *Methodol.* **71**: 393–424.
- 32 Sun D., Xi Y., Rodriguez B., Park H. J., Tong P., Meong M., Goodell M. A., Li W., 2014 MOABS: model based
33 analysis of bisulfite sequencing data. *Genome Biol.* **15**: R38.
- 34 Sun S., Yu X., 2016 HMM-Fisher: identifying differential methylation using a hidden Markov model and
35 Fisher’s exact test. *Stat. Appl. Genet. Mol. Biol.* **15**: 55–67.
- 36 Tierney L., Kadane J. B., 2012 Accurate Approximations for Posterior Moments and Marginal Densities. *J.*
37 *Am. Stat. Assoc.*
- 38 Wang J., Duncan D., Shi Z., Zhang B., 2013 WEB-based GENE SeT Analysis Toolkit (WebGestalt): update
39 2013. *Nucleic Acids Res.* **41**: W77–83.
- 40 Wu T., Ye Y., Min S.-Y., Zhu J., Khobahy E., Zhou J., Yan M., Hemachandran S., Pathak S., Zhou X. J., Andreeff

1 M., Mohan C., 2014 Prevention of murine lupus nephritis by targeting multiple signaling axes and
2 oxidative stress using a synthetic triterpenoid. *Arthritis Rheumatol.* (Hoboken, N.J.) **66**: 3129–39.

3 Wu H., Xu T., Feng H., Chen L., Li B., Yao B., Qin Z., Jin P., Conneely K. N., 2015 Detection of differentially
4 methylated regions from whole-genome bisulfite sequencing data without replicates. *Nucleic Acids*
5 *Res.* **43**: gkv715.

6 Wu H., Zhao M., Tan L., Lu Q., 2016 The key culprit in the pathogenesis of systemic lupus erythematosus:
7 Aberrant DNA methylation. *Autoimmun. Rev.* **15**: 684–689.

8 Yu X., Sun S., 2016a HMM-DM: identifying differentially methylated regions using a hidden Markov model.
9 *Stat. Appl. Genet. Mol. Biol.* **15**: 69–81.

10 Yu X., Sun S., 2016b Comparing five statistical methods of differential methylation identification using
11 bisulfite sequencing data. *Stat. Appl. Genet. Mol. Biol.* **15**: 173–91.

12 Zhang W., Spector T. D., Deloukas P., Bell J. T., Engelhardt B. E., 2015 Predicting genome-wide DNA
13 methylation using methylation marks, genomic position, and DNA regulatory elements. *Genome Biol.*
14 **16**: 14.

15 Ziller M. J., Müller F., Liao J., Zhang Y., Gu H., Bock C., Boyle P., Epstein C. B., Bernstein B. E., Lengauer T.,
16 Gnirke A., Meissner A., 2011 Genomic distribution and inter-sample variation of non-CpG
17 methylation across human cell types. (D Schübeler, Ed.). *PLoS Genet.* **7**: e1002389.

18 Ziller M. J., Hansen K. D., Meissner A., Aryee M. J., 2015 Coverage recommendations for methylation
19 analysis by whole-genome bisulfite sequencing. *Nat. Methods* **12**: 230–2, 1 p following 232.

20 Zoghbi H. Y., Beaudet A. L., 2016 Epigenetics and Human Disease. *Cold Spring Harb. Perspect. Biol.* **8**:
21 a019497.

22